

AniPortraitGAN: Animatable 3D Portrait Generation from 2D Image Collections

Yue Wu^{1*} Sicheng Xu^{2*} Jianfeng Xiang^{3,2} Fangyun Wei²
 Qifeng Chen¹ Jiaolong Yang^{2†} Xin Tong²
¹HKUST ²Microsoft Research Asia ³Tsinghua University



Figure 1: Our method is a new 3D-aware GAN that can generate diverse *virtual* human portraits (512×512) with explicitly controllable 3D camera viewpoints, facial expression, head pose, and shoulder movements. It is trained on unstructured 2D images without any 3D or video data. (Best viewed with zoom; see our [project page](#) for videos of more samples)

Abstract

Previous animatable 3D-aware GANs for human generation have primarily focused on either the human head or full body. However, head-only videos are relatively uncommon in real life, and full body generation typically does not deal with facial expression control and still has challenges in generating high-quality results. Towards applicable video avatars, we present an animatable 3D-aware GAN that generates portrait images with controllable facial expression, head pose, and shoulder movements. It is a generative model trained on unstructured 2D image collections without using 3D or video data. For the new task, we base our method on the generative radiance manifold representation and equip it with learnable facial and head-shoulder deformations. A dual-camera rendering and adversarial learning scheme is proposed to improve the quality of the generated faces, which is critical for portrait images. A pose deformation processing network is developed

to generate plausible deformations for challenging regions such as long hair. Experiments show that our method, trained on unstructured 2D images, can generate diverse and high-quality 3D portraits with desired control over different properties.

1. Introduction

The automatic creation of animatable 3D human characters has become an increasingly important topic with a range of applications including video conferencing, movie production, and gaming. The related techniques have undergone a significant growth recently, with a variety of promising methods being proposed [1, 7, 8, 16, 21, 38, 40, 42, 51, 52, 54, 60, 62].

Among these techniques, 3D-aware generation methods have emerged as a particularly promising avenue [1, 12, 21, 38, 51, 52, 60, 62]. These methods can leverage the ample availability of unstructured 2D data for 3D human generative learning without the need for 3D scans or multiview human images which are difficult to acquire at scale. Typ-

*Equal contribution. Work done when YW was an intern at MSRA.

†Corresponding author and project lead.

ically, these methods employ Generative Adversarial Networks (GANs) [18] for unsupervised training, use neural implicit fields [34] as the 3D representation, and incorporate priors from 3D face and body parametric models [4, 29, 32] for character control.

Despite their promising potential, existing animatable 3D-aware human generation focuses on either the human head or full body and have encountered limitations in their applicability. Head generation methods [51, 52, 60, 62] can produce high-quality face and hair with controllable facial expression. Unfortunately, videos featuring only a human head are relatively uncommon in everyday life, and therefore these methods are less applicable in practical scenarios. Full body generation [1, 12, 21, 38, 65], on the other hand, also generates torso and limbs with explicit pose control. However, generating high-quality full body human is still challenging due to the complexity of body motion. In addition, the facial region is often underrepresented in these full body methods and there is no expression control.

Our paper presents a new 3D-aware generation method that is the first to focus on animatable generation of the human head and shoulder regions. Our method enables fine-grained control over facial expressions as well as head and shoulder movements, making it well suited for real-world applications such as video conferencing and virtual presenters. Like previous 3D-aware GANs, our method is trained on unstructured 2D image sets.

For this new task, we follow previous methods to train neural radiance generation with 3D parametric model priors in a GAN training scheme. We base our method on the 3D-aware GAN framework of GRAM [9, 61] and follow AniFaceGAN [60] for facial expression control using 3D morphable model (3DMMs) priors. For head and shoulder control, we incorporate the SMPL [32] body model for deformation guidance. However, we found that naively extending these existing techniques to our animatable head-shoulder portrait generation task is deficient. One prominent issue is about face quality, which is of paramount importance for visual communications. To handle the complex image distribution caused by the large variations of head position and orientation, we propose a dual-camera rendering and adversarial learning scheme for training. An additional dynamic camera is placed around human head pointing at head center to render faces for discrimination, which significantly improves the face generation quality. Another issue is the SMPL-guided human body deformation, for which we identified that the commonly-used strategy based on linear blending skinning failed to generate convincing results for human characters with long hair. Sharp discontinuities will occur under head rotation for hair region, leading to significant artifacts. To tackle this issue, we propose a pose deformation volume processing module to learn better deformations, which stabilizes GAN training and produces vi-

sually plausible results.

We train our model on a head-shoulder portrait dataset called SHHQ-HS, which is constructed by cropping and superresolving the 40K human body images in the SHHQ dataset [15]. We show that our method can generate diverse and high-quality 3D portrait images with flexible control of different properties including facial expressions and head-shoulder poses.

Our contributions can be summarized as follows:

- We propose the first animatable 3D-aware portrait GAN that generates head and shoulder regions with facial expression and head-shoulder motion control. We believe generating such animatable human characters is a missing piece of 3D-aware human GANs for real-world applications like video conferencing and virtual presenters.
- We propose a dual-camera rendering and adversarial learning scheme that gives rise to high-quality face generation comparable to previous head-only 3D-aware GANs.
- We propose a pose deformation processing module which achieves smooth and plausible pose-driven deformation for human hair.

2. Related Work

3D-aware Image Generation 3D-aware image generative models aim to generate images that allow for the explicit control of 3D camera viewpoint, training only on 2D images. Most existing works are based on the GAN framework for generative modeling. Early approaches to this problem use 3D convolutions to generate 3D feature volumes and project them to 2D plane for image generation [35, 36]. Recently, methods based on more explicit 3D representations and differentiable rendering have become popular [5, 6, 9, 11, 19, 30, 37, 39, 44, 47, 49, 50, 53, 61, 66]. The widely-used 3D representations are NeRF [34] and its variants, for their strong capability to model real-world 3D scenes.

Among these NeRF-based GANs, some use volume rendering to directly generate the final images, which ensures strong 3D consistency among different views but often has high computation cost. Others apply 2D convolutions on the rendered low-resolution 2D images or feature maps for up-sampling, which significantly reduces the computation cost but sacrifices multiview consistency for the generated instances. Our method uses the high-resolution radiance manifold representation of [61], which can generate high resolution images with strong multiview consistency.

Controllable Human Head and Body Generation

Adding explicit controls to face and body generative mod-

eling has received much attention in recent years. Existing works have primarily focused on either the head [1, 8, 51, 52, 60, 62] or the whole body [1, 7, 12, 21, 38], with priors from 3D parametric models being commonly incorporated to achieve semantically meaningful control. For expression control, most head GANs [1, 8, 51, 52, 60] often incorporate 3D morphable models (3DMMs) [4] or FLAME models [29] in their training process. Whole-body GANs typically rely on the SMPL model [32] for body pose animation (with a few exceptions such as [38] that use body skeleton), and they often do not address facial expression control. This work deals with a new human generation task: generating portrait figures that contain head and shoulder, with controllable facial expression and head-shoulder poses.

Human Image and Video Manipulation Our method is also related to human image and video manipulation approaches [13, 16, 17, 23, 27, 31, 42, 43, 48, 55, 58, 59, 63, 64] that also produce human animation videos. However, the goal and underlying techniques of these methods differ significantly from ours. These methods aim to animate the human character in the given image or video, and are typically trained in a supervised manner using videos or image pairs. In contrast, we deal with human generative modeling and novel character creation, training on unstructured still images in an unsupervised or weakly-supervised fashion.

3. Method

Our goal is to generate human portrait images containing human head and shoulder regions by training on a given 2D image collection. As in a standard GAN setup, we sample random latent codes and map them to the final output image. The input to our generator consists of multiple latent codes, which corresponds to different properties of the generated human, and the camera viewpoint. The output is a human portrait image carrying the desired properties.

Figure 2 presents an overview of our method. The overall pipeline follows the popular paradigm of canonical neural radiance representation in combination with (inverse) deformation [1, 21, 60].

3.1. Latent Codes

Our latent codes contain an identity code $\mathbf{z}_{id} \in \mathbb{R}^{d_i}$ for human shape, an expression code $\mathbf{z}_{exp} \in \mathbb{R}^{d_e}$ for facial expression, an $\mathbf{z}_{pose} \in \mathbb{R}^{d_p}$ for head and shoulder pose, and an additional noise $\boldsymbol{\varepsilon} \in \mathbb{R}^{d_\varepsilon}$ controlling other attributes such as appearance. To achieve semantically meaningful control, we incorporate priors 3D human parametric models and align our latent space with theirs. Specifically, our identity code \mathbf{z}_{id} is designed as the concatenation of 3DMM [41] face identity coefficient and SMPL [32] body shape coefficient. The pose code \mathbf{z}_{pose} is a reduced SMPL pose parameter, which consists of the joint transformations of 6 joints:

head, neck, left and right collars, and left and right shoulders. The expression code \mathbf{z}_{exp} is the same as 3DMM expression coefficient.

3.2. Canonical Radiance Manifolds

Our method utilizes the radiance manifolds [9, 61] to represent canonical humans. This representation regulates radiance field learning and rendering on a set of learned implicit surfaces in the 3D volume. It can generate high-quality human faces with strict multiview consistency. With manifold superresolution [61], it can generate high-resolution images efficiently without sacrificing multiview consistency.

Concretely, we apply three networks for radiance generation. A manifold prediction MLP \mathcal{M} takes a point \mathbf{x} in the canonical space as input and predicts a scalar s :

$$\mathcal{M} : \mathbf{x} \in \mathbb{R}^3 \rightarrow s \in \mathbb{R}. \quad (1)$$

It models a scalar field that defines the surfaces. A radiance generation MLP ϕ generates the color and opacity for points on the surfaces given the identity codes \mathbf{z}_{id} , noise $\boldsymbol{\varepsilon}$ and view direction \mathbf{d} :

$$\phi : (\mathbf{x}, \mathbf{z}_{id}, \boldsymbol{\varepsilon}, \mathbf{d}) \in \mathbb{R}^{d_i+d_\varepsilon+6} \rightarrow (\mathbf{c}, \alpha) \in \mathbb{R}^4. \quad (2)$$

A manifold superresolution CNN \mathcal{U} upsamples the flattened and discretized radiance maps \mathbf{R}_{lr} to high-resolution ones \mathbf{R}_{hr} :

$$\mathcal{U} : (\mathbf{z}_{id}, \boldsymbol{\varepsilon}, \mathbf{R}_{lr}) \rightarrow \mathbf{R}_{hr}, \quad (3)$$

for which we use a $128^2 \rightarrow 512^2$ upsampling setting in this paper. For more technical details, we refer the readers to [9, 61].

3.3. Deformation Fields

For each sampled 3D point in the target space with desired head-shoulder pose and facial expression, we apply deformations to transform them to the canonical space for radiance retrieval. A two-stage deformation scheme is used to neutralize pose and expression.

3.3.1 Pose Deformation Generator

We incorporate the SMPL model [32] and use its linear blend skinning (LBS) scheme [28] to guide our deformation. Given the shape code \mathbf{z}_{id} and pose code \mathbf{z}_{pose} , a posed human body mesh can be constructed using SMPL. The SMPL model provides a pre-defined skinning weight vector $\mathbf{w} \in \mathbb{R}^{N_J}$ for each vertex on the body surface, where N_J is the joint number.

A simple approach for propagating body surface deformation to the full 3D space is assigning any point the skinning weights of its closest body surface vertex and use them

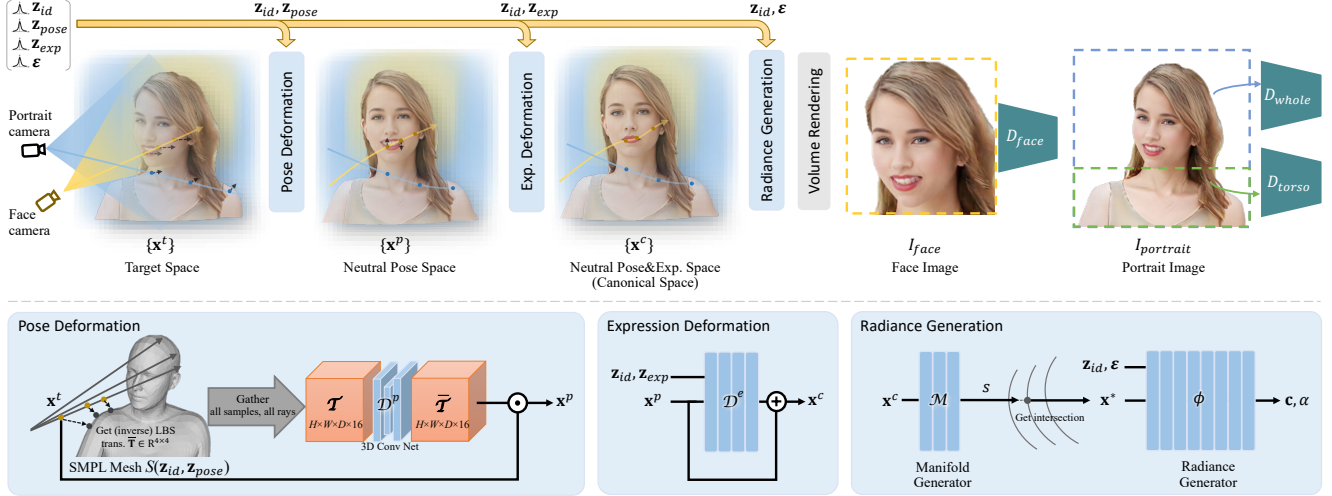


Figure 2: Method overview. *Top*: the pipeline of our controllable 3D-aware portrait GAN. For training, we apply a dual-camera rendering scheme with two images separated rendered and three discriminators employed. *Bottom*: structures of the deformation and radiance generation modules (the radiance manifold super-resolution step is omitted for simplicity; see text for details).

to deform it. In fact, this strategy is widely used in state-of-the-art animatable human body modeling and generation methods [2, 12, 17, 22, 42, 65]. While it yields plausible results for existing full-body synthesis method, we find it incurs significant visual defects in high-resolution portrait synthesis. For human characters with long hairs, this native strategy leads to sharp deformation discontinuity for the hair regions above shoulders (see Fig. 6).

We propose a deformation volume processing module to tackle this issue. Specifically, for a point \mathbf{x}^t in the target space with the skinning weight vector \mathbf{w} retrieved from the closest SMPL body vertex, the deformed point can be calculated using inverse LBS:

$$\mathbf{x}^p = LBS^{-1}(\mathbf{x}^t, \mathbf{w}) = \bar{\mathbf{T}} \cdot \mathbf{x}^t = (\sum_{j=1}^{N_J} w_j \mathbf{T}_j) \cdot \mathbf{x}^t, \quad (4)$$

where $\bar{\mathbf{T}} \in \mathbb{R}^{4 \times 4}$ is a transformation matrix computed by linearly blending the SMPL joint transformations $\mathbf{T}_j \in SE(3)$. We gather the reshaped transformation matrices for the sampled points of all $H \times W$ rays into a tensor $\mathcal{T} \in \mathbb{R}^{H \times W \times D \times 16}$, where D is the number of sampled points per ray, and apply a 3D CNN \mathcal{D}^p to process it:

$$\mathcal{D}^p : \mathcal{T} \in \mathbb{R}^{H \times W \times D \times 16} \rightarrow \bar{\mathcal{T}} \in \mathbb{R}^{H \times W \times D \times 16}. \quad (5)$$

After processing, we reshape the transformations back and apply them to the sampled points to accomplish pose deformation.

3.3.2 Expression Deformation Generator

We further apply a deformation to neutralize the facial expression from the target expression. Following [60], we introduce a deformation field which is guided by the 3DMM

model [41]. Specifically, an MLP \mathcal{D}^e is employed to deform the points in the pose-aligned space:

$$\mathcal{D}^e : (\mathbf{x}^p, \mathbf{z}_{id}, \mathbf{z}_{exp}) \rightarrow \mathbf{x}^c. \quad (6)$$

This deformation network will be trained to generate faces with expressions following 3DMM:

$$\mathbf{S} = \mathbf{S}(\mathbf{z}_{id}, \mathbf{z}_{exp}) = \bar{\mathbf{S}} + \mathbf{B}_{id} \mathbf{z}_{id} + \mathbf{B}_{exp} \mathbf{z}_{exp}, \quad (7)$$

where $\bar{\mathbf{S}}$ is the 3DMM mean face, and \mathbf{B}_{id} and \mathbf{B}_{exp} are the PCA basis for identity and expression, respectively. Training details can be found in later sections.

In total, our generator G has five sub-nets: \mathcal{M} , ϕ , \mathcal{U} , \mathcal{D}^p , and \mathcal{D}^e . For final image rendering, we calculate M intersection points $\{\mathbf{x}_i^*\}$ between a deformed ray \mathbf{r} (point samples) and the canonical manifolds. We then obtain the color and occupancy of $\{\mathbf{x}_i^*\}$ by sampling the radiance map \mathbf{R}_{hr} , and composite the color via:

$$C(\mathbf{r}) = \sum_{i=1}^M T(\mathbf{x}_i^*) \alpha(\mathbf{x}_i^*) \mathbf{c}(\mathbf{x}_i^*), \quad T(\mathbf{x}_i^*) = \prod_{k < i} (1 - \alpha(\mathbf{x}_k^*)). \quad (8)$$

3.4. Dual-Camera Discriminators

Previous 3D-aware head GANs have demonstrated striking face generation quality by carefully center-aligning the generated and real face images for training. In our case, however, the head region constitutes part of the portrait image and its spatial position and orientation vary significantly. Simply applying a whole-image discriminator cannot offer adequate supervision for high-quality face generation, which is crucial for portrait images.

A straightforward remedy is to crop and align the faces in the rendered images and apply a local face discriminator. However, since image resampling operators are inherently low-pass, such an image-space cropping strategy introduces blur to the cropped faces, which is detrimental to GAN training especially for a multi-discriminator setup. In this work, we design a dual-camera rendering scheme for GAN training. In addition to the main camera for full portrait image rendering, we add another camera for face rendering, which is placed around human head pointing at head center, as shown in Fig. 2. It is designed to have the same local coordinate system as in previous 3D-aware head GANs [9, 60], and its position can be readily computed using the deformed SMPL head. Another possible idea is blending the output of two separate generators for face and body, as in some 2D human generation methods [14]. But applying this strategy to the 3D, animatable case seems not straightforward.

Adding a dedicated face camera for training not only avoids image resampling and provides a more direct supervision to the canonical radiance manifolds, but also enables higher-resolution face rendering for adversarial learning. Consequently, the radiance generator receives stronger supervision for the face region. We apply two image discriminators D_{whole} and D_{face} for the rendered portrait image and face image with these two cameras, respectively. We also add another local discriminator D_{torso} , which takes the lower 1/4 part of the rendered portrait images as input.

3.5. Training Losses

Adversarial Learning. We apply the non-saturating GAN loss with R1 regularization [33] for the 3D-aware image generator and all three discriminators D_{whole} , D_{face} , and D_{torso} . We empirically set the balancing weights to be $\lambda_{whole} = 0.1$, $\lambda_{face} = 1.0$ and $D_{torso} = 0.5$, respectively.

Deformation Learning Following [60], we use a 3D landmark loss and imitation loss to gain expression control with 3DMM guidance. The landmark loss enforces the generated face image to have similar 3D facial landmarks to the 3DMM face constructed with the input identity and expression codes:

$$\mathcal{L}_{lm} = \|f_{lm}(\mathbf{S}(\mathbf{z}_{id}, \mathbf{z}_{exp})), f_{lm}(\mathbf{S}(\hat{\mathbf{z}}_{id}, \hat{\mathbf{z}}_{exp}))\|_2^2, \quad (9)$$

where $\hat{\mathbf{z}}_{id}$, $\hat{\mathbf{z}}_{exp}$ are the 3DMM coefficients estimated from the generated image using a face reconstruction network [10] and f_{lm} denotes a simple facial landmark extraction function. For deformation imitation, we enforce the displacement of an input point \mathbf{x}^p to follow its nearest point \mathbf{x}_{ref}^p on the 3DMM mesh:

$$\mathcal{L}_{3DMM} = \left\| \left(\mathcal{D}^e(\mathbf{x}^p) - \mathbf{x}^p \right) - \Delta \mathbf{x}_{ref}^p \right\|_2^2, \quad (10)$$

where $\Delta \mathbf{x}_{ref}^p = -\mathbf{B}_{exp} \mathbf{z}_{exp}(\mathbf{x}_{ref}^p)$ is the 3DMM-derived (inverse) deformation of point \mathbf{x}_{ref}^p .

We impose several regularizations on the deformations. First, we encourage the deformations to be as smooth as possible. For pose deformation processing, we apply

$$\mathcal{L}_{pose_smooth} = \|\mathcal{D}^p(\mathcal{T}) - AvgPool(\mathcal{T})\|_2^2. \quad (11)$$

where $AvgPool$ is a fixed average pooling operator. For the expression deformation, we apply

$$\mathcal{L}_{exp_smooth} = \|\mathcal{D}^e(\mathbf{x}^p) - \mathcal{D}^e(\mathbf{x}^p + \Delta)\|_2^2, \quad (12)$$

where Δ is a small random perturbation. Finally, a minimal deformation constraint is applied for the expression deformation:

$$\mathcal{L}_{exp_minimal} = \|\mathbf{x}^p - \mathcal{D}^e(\mathbf{x}^p)\|_2^2. \quad (13)$$

3.6. Training Strategy

We employ a two-stage training strategy to train our model. At the first stage, we train a low-resolution image generator and the corresponding discriminators. Both the face and portrait branches generate 128×128 images. All sub-networks are trained except for the manifold super-resolution CNN \mathcal{U} . For the second stage, we generate 512×512 portrait images and 256×256 faces. We randomly initialize and train \mathcal{U} as well as the high-resolution discriminators with all other sub-networks frozen.

4. Experiments

Training Data We build a training set by processing the human images in the SHHQ dataset [15]. SHHQ contains 40K full-body images of 1024×512 resolution. To obtain high-quality head-shoulder portraits, we first fit SMPL models on the SHHQ images using the method of [24]. Then, we crop the images and align them using the projected head and neck joints. The cropped portrait images are about 256×256 resolution. We upsample them using the super-resolution methods of [56] and [57] to 1024×1024 , followed by downsampling them to 512×512 . Finally, the backgrounds are removed by applying the provided segmentation masks. We call this dataset *SHHQ-HS*.

Implementation Details Our manifold predictor and radiance generator follow the implementations of [9]. 24 radiance manifolds are used as in [9]. The manifold super-resolution net \mathcal{U} is a smaller CNN compared to that in [61]. The pose deformation CNN \mathcal{D}^p has two 3D conv layers with kernel size $9 \times 9 \times 5$. The expression deformation network \mathcal{D}^e is the same as [60]. See Fig. 10 for more details. For all experiments, we use the Adam optimizer [26] for training. Prior to training, we estimate the identity, pose and expression coefficients as well as camera poses for the images in



Figure 3: Results with controllable camera viewpoint, facial expression, head pose and shoulder pose. (Best viewed with zoom)

the dataset using 3DMM and SMPL fitting [10, 24]. During training, we randomly sample latent codes \mathbf{z}_{id} , \mathbf{z}_{pose} , and \mathbf{z}_{exp} and camera pose θ from the estimated distributions, and sample ϵ from a normal distribution.

Runtime. With an unoptimized implementation, our method takes about 0.87 seconds to generate one 512×512 image from a set of given latent codes, evaluated on a

NVIDIA RTX A6000 GPU.

4.1. Generation Results

Figure 1 and 9 present some generated portraits from our method trained on SHHQ-HS. The results are diverse and of high-quality, with camera viewpoint, facial expression, head rotation, and shoulder pose explicitly controlled. Figure 3 shows the generated results for which we control one

Table 1: Quantitative comparison with state-of-the-art 3D-aware GANs on SHHQ-HS. Note that EG3D and GRAM-HD do not provide any expression and pose control, and AniFaceGAN only generates head images. FID and KID ($\times 100$) are computed with 20K randomly generated images and 20K real ones. *: Although EG3D has lowest FID and KID scores, it often generates planar geometry; see Fig. 4. †: AniFaceGAN is trained on 128^2 by [60] and evaluated with 256^2 rendering.

Method	Face 256^2		Full 512^2	
	FID↓	KID↓	FID↓	KID↓
EG3D	5.63*	0.20*	6.81*	0.26*
GRAM-HD _{64→512}	8.01	0.41	7.75	0.29
GRAM-HD _{128→512}	8.14	0.30	8.82	0.28
AniFaceGAN	11.56†	0.66†	N/A	N/A
Ours	7.64	0.43	10.10	0.43

property out of the four while randomly changing the others. Our method achieves consistent control for all the four properties for different identities. More results can be found in the suppl. video.

4.2. Comparison with Previous Methods

We compare our method with three state-of-the-art 3D-aware GANs: EG3D [5], GRAM-HD [61] and AniFaceGAN [60]. Note that to our knowledge, there is no previous work that deals with the animatable head-shoulder portrait generation task in this paper, and hence we compare with these three methods for reference purpose only.

Table 1 shows the FID [20] and KID [3] metrics evaluated on both the full portrait images and the face regions. Our method has comparable scores with EG3D and GRAM-HD for face and slightly lower scores on full image. Note that although EG3D has lowest scores, we found that it often generates poor geometry: the portrait surfaces are sometimes nearly planar and the visual parallax is wrong when changing viewing angles (Fig. 4). Visually inspected, our image quality is comparable with EG3D and GRAM-HD and the portraits have correct geometry, as shown in Fig. 4. Figure 5 compares the results from AniFaceGAN and our method. Clearly, our method can generate and control much larger region.

4.3. Ablation Study

We then conduct ablation studies to validate the effectiveness of our algorithm design. For efficiency, we quantitatively evaluate the generation quality on the 128^2 resolution (i.e., without manifold super-resolution) and compute the FID and KID metrics using 5K generated and real images. The results are summarized in Table 2.

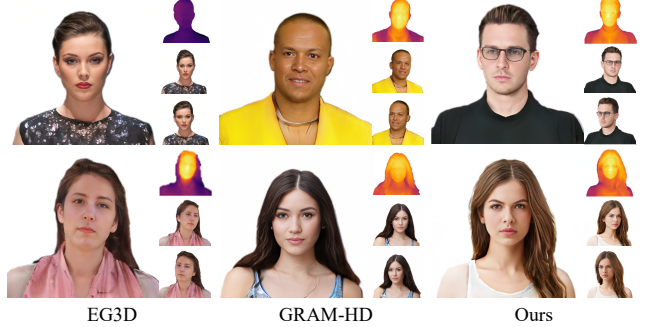


Figure 4: Visual comparison with state-of-the-art 3D-aware GANs on SHHQ-HS. Our results have similar visual quality to existing 3D-aware GANs that do not handle expression and pose control. (Best viewed with zoom)



Figure 5: Visual comparison with AniFaceGAN.

Discriminator Combinations Our full method uses three discriminators for training: D_{whole} , D_{face} , and D_{torso} . Table 2 shows that using D_{whole} alone is clearly deficient, as demonstrated by the poor FID and KID scores for face. Combining D_{whole} with D_{face} (i.e., without D_{torso}) significantly improved the face quality, but the full portrait quality gets much worse than using only D_{whole} . The FID and KID scores of our full method are low for both the full portrait images and face regions.

Dual-Camera Training To validate our dual-camera rendering and adversarial learning scheme, we train a variant of our method without a separate face camera for training. For this variant, we locate faces in the rendered head-shoulder portrait images during training, and apply face discriminator on cropped and aligned faces. As we can see from Table 2, our method with dual cameras for training significantly outperforms such an image cropping strategy in terms of face quality with slightly lower full-image FID. Some visual examples can be found in Fig. 6.

Pose Deformation Processing Module Table 2 (last row) shows that removing our pose deformation processing CNN \mathcal{D}^p , which degrades to a simple skinning weight assignment strategy, leads to a quality drop for the full portrait image. Figure 6 visually compares two typical generation results. Without \mathcal{D}^p , sharp discontinuities will occur for long hairs

Table 2: Ablation study on discriminator settings and the pose deformation processing CNN \mathcal{D}^p . *: w/o dual-camera training, i.e., applying D_{face} on faces cropped from the rendered portrait images. FID and KID ($\times 100$) are computed with 5K randomly generated images and 5K real ones.

Method				Face 128 ²		Full 128 ²	
D_{whole}	D_{face}	D_{torso}	\mathcal{D}^p	FID↓	KID↓	FID↓	KID↓
✓	✓	✓	✓	11.26	0.57	16.68	0.97
✓	×	×	✓	23.00	1.48	18.11	1.06
✓	✓	×	✓	10.68	0.52	22.58	1.49
✓	*	✓	✓	17.89	1.33	14.05	0.69
✓	✓	✓	×	10.88	0.52	19.27	1.23

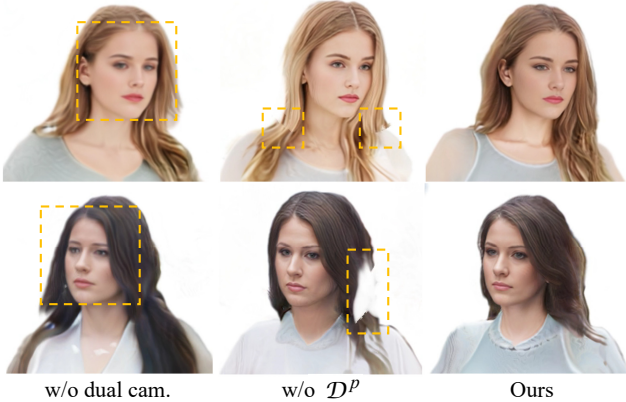


Figure 6: Ablation on the dual camera training scheme and pose deformation CNN \mathcal{D}^p (trained on 128² and rendered on 256²). See the suppl. video for more results.

under head movements, whereas our full method produces plausible results without obvious artifacts.

4.4. Talk Video Generation

We further test our trained model on the task of generating videos of talking portraits driven by real videos. Specifically, we selected some talk videos from the 300-VW dataset [45] and track the 3DMM expression and SMPL head-shoulder pose using the methods of [10] and [24], respectively. Simple temporal smoothing is applied on the estimation results. Then we transfer the tracked results to our generated human characters to obtain virtual talk videos. Figure 7 shows some typical examples of our results where a generated virtual character moves following the real person. See the supplementary video for continuous animations of our results.

5. Conclusion

We have presented a novel 3D-aware GAN for animatable head-shoulder portrait generation, a new task not addressed by previous methods. We identified several key is-

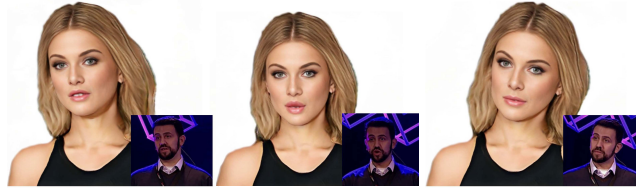


Figure 7: Talk video generation driven by real person.

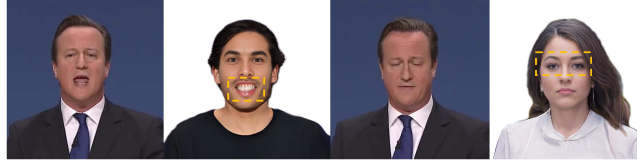


Figure 8: Limitations on more extreme expressions and closed eyes. For each image pair, the left one is the reference image and the right one is the animated result.

sues when extending existing techniques to this new task and proposed targeted algorithms to tackle them. We demonstrate that by training a corpus of unstructured 2D images, our method can generate diverse and high-quality 3D portraits with controllable facial expression as well as head and shoulder movements. We believe our work represents one step forward towards auto-creating video avatars for real-world applications.

Limitations Our method still has several limitations. It may produce artifacts under human poses and expressions that are not present in training data distribution, as shown in Fig. 8. In fact, the facial expression variation in SHHQ-HS is rather limited, lacking images with extreme expression and closed eyes. The visual quality of the inner mouth region (e.g., teeth) is not satisfactory, which is also partially due to limited data samples. Additionally, our current method lacks the ability to control other attributes, such as eye gaze and environment lighting. We plan to further explore and address these issues in our future research.

Ethics and responsible AI considerations This work aims to design an animatable 3D-aware human portrait generation method for the application of virtual avatars. It is not intended to create content that is used to mislead or deceive. However, it could still potentially be misused. We condemn any behaviors of creating misleading or harmful contents and are interested in applying this technology for advanced forgery detection. Currently, the images generated by this method contain visual artifacts that can be easily identified. The method's performance is affected by the biases in the training data. One should be careful about the data collection process and ensure unbiased distributions of race, gender, age, among others.

References

- [1] Alexander Bergman, Petr Kellnhofer, Wang Yifan, Eric Chan, David Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. In *Advances in Neural Information Processing Systems*, 2022. [1](#), [2](#), [3](#)
- [2] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *Advances in Neural Information Processing Systems*, pages 12909–12922, 2020. [4](#)
- [3] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018. [7](#)
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Annual Conference on Computer Graphics and Interactive Techniques*, pages 187–194, 1999. [2](#), [3](#)
- [5] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. [2](#), [7](#)
- [6] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. [2](#), [13](#)
- [7] Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J Black, Andreas Geiger, and Otmar Hilliges. gDNA: Towards generative detailed neural avatars. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20427–20437, 2022. [1](#), [3](#)
- [8] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5154–5163, 2020. [1](#), [3](#)
- [9] 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 10673–10683, 2022. [2](#), [3](#), [5](#)
- [10] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019. [5](#), [6](#), [8](#)
- [11] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W Taylor, and Joshua M Susskind. Unconstrained scene generation with locally conditioned radiance fields. In *IEEE/CVF International Conference on Computer Vision*, 2021. [2](#)
- [12] Zijian Dong, Xu Chen, Jinlong Yang, Michael J Black, Otmar Hilliges, and Andreas Geiger. AG3D: Learning to generate 3d avatars from 2d image collections. *arXiv preprint arXiv:2305.02312*, 2023. [1](#), [2](#), [3](#), [4](#)
- [13] Michail Christos Doukas, Stefanos Zafeiriou, and Viktoriia Sharmanska. Headgan: One-shot neural head synthesis and editing. In *IEEE/CVF International Conference on Computer Vision*, pages 14398–14407, 2021. [3](#)
- [14] Anna Frühstück, Krishna Kumar Singh, Eli Shechtman, Niloy J Mitra, Peter Wonka, and Jingwan Lu. Insetgan for full-body image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7723–7732, 2022. [5](#)
- [15] Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen Change Loy, Wayne Wu, and Ziwei Liu. Stylegan-human: A data-centric odyssey of human generation. In *European Conference on Computer Vision*, pages 1–19, 2022. [2](#), [5](#)
- [16] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. [1](#), [3](#)
- [17] Xiangjun Gao, Jiaolong Yang, Jongyoo Kim, Sida Peng, Zicheng Liu, and Xin Tong. MPS-NeRF: Generalizable 3d human rendering from multiview images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. [3](#), [4](#)
- [18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, 2014. [2](#)
- [19] 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. [2](#)
- [20] 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. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017. [7](#)
- [21] Fangzhou Hong, Zhaoxi Chen, Yushi Lan, Liang Pan, and Ziwei Liu. EVA3D: Compositional 3d human generation from 2d image collections. In *International Conference on Learning Representations*, 2023. [1](#), [2](#), [3](#)
- [22] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3093–3102, 2020. [4](#)
- [23] Kaiwen Jiang, Shu-Yu Chen, Feng-Lin Liu, Hongbo Fu, and Lin Gao. Nerffacediting: Disentangled face editing in neural radiance fields. In *ACM SIGGRAPH Asia*, pages 1–9, 2022. [3](#)
- [24] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human model fitting towards in-the-wild 3d human pose estimation. In *International Conference on 3D Vision*, pages 42–52, 2021. [5](#), [6](#), [8](#)
- [25] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. [13](#)
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. [5](#)

- [27] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances in Neural Information Processing Systems*, 34:24741–24752, 2021. [3](#)
- [28] John P Lewis, Matt Cordner, and Nickson Fong. Pose space deformation: a unified approach to shape interpolation and skeleton-driven deformation. In *Annual Conference on Computer Graphics and Interactive Techniques*, pages 165–172, 2000. [3](#)
- [29] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Transactions on Graphics*, 36(6):194–1, 2017. [2](#), [3](#)
- [30] Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. Towards unsupervised learning of generative models for 3d controllable image synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5871–5880, 2020. [2](#)
- [31] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Transactions on Graphics*, 40(6):1–16, 2021. [3](#)
- [32] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM transactions on Graphics*, 34(6):1–16, 2015. [2](#), [3](#)
- [33] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International Conference on Machine Learning*, pages 3481–3490, 2018. [5](#)
- [34] 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 *European Conference on Computer Vision*, pages 405–421, 2020. [2](#)
- [35] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3D representations from natural images. In *IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. [2](#)
- [36] Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy Mitra. BlockGAN: Learning 3d object-aware scene representations from unlabelled images. In *Advances in Neural Information Processing Systems*, 2020. [2](#)
- [37] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. [2](#)
- [38] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Unsupervised learning of efficient geometry-aware neural articulated representations. In *European Conference on Computer Vision*, pages 597–614, 2022. [1](#), [2](#), [3](#)
- [39] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. [2](#)
- [40] Hao Ouyang, Bo Zhang, Pan Zhang, Hao Yang, Jiaolong Yang, Dong Chen, Qifeng Chen, and Fang Wen. Real-time neural character rendering with pose-guided multiplane images. In *European Conference on Computer Vision*, pages 192–209, 2022. [1](#)
- [41] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *IEEE International Conference on Advanced Video and Signal based Surveillance*, pages 296–301, 2009. [3](#), [4](#)
- [42] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *IEEE/CVF International Conference on Computer Vision*, pages 14314–14323, 2021. [1](#), [3](#), [4](#)
- [43] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *IEEE/CVF International Conference on Computer Vision*, pages 13759–13768, 2021. [3](#)
- [44] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems*, 2020. [2](#)
- [45] Jie Shen, Stefanos Zafeiriou, Grigoris G Chrysos, Jean Kossai, Georgios Tzimiropoulos, and Maja Pantic. The first facial landmark tracking in-the-wild challenge: Benchmark and results. In *IEEE International Conference on Computer Vision Workshops*, pages 50–58, 2015. [8](#)
- [46] Wenzhe Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016. [13](#)
- [47] Yichun Shi, Divyansh Aggarwal, and Anil K Jain. Lifting 2d stylegan for 3d-aware face generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6258–6266, 2021. [2](#)
- [48] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *Advances in Neural Information Processing Systems*, 32, 2019. [3](#)
- [49] Ivan Skorokhodov, Aliaksandr Siarohin, Yinghao Xu, Jian Ren, Hsin-Ying Lee, Peter Wonka, and Sergey Tulyakov. 3d generation on imagenet. In *International Conference on Learning Representations*, 2023. [2](#)
- [50] Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. Epigraf: Rethinking training of 3d gans. In *Advances in Neural Information Processing Systems*, 2022. [2](#)
- [51] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. [1](#), [2](#), [3](#)
- [52] Keqiang Sun, Shangzhe Wu, Zhaoyang Huang, Ning Zhang, Quan Wang, and HongSheng Li. Controllable 3d face synthesis with conditional generative occupancy fields. In *Advances in Neural Information Processing Systems*, 2022. [1](#), [2](#), [3](#)

- [53] Ayush Tewari, Xingang Pan, Ohad Fried, Maneesh Agrawala, Christian Theobalt, et al. Disentangled3d: Learning a 3d generative model with disentangled geometry and appearance from monocular images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1516–1525, 2022. [2](#)
- [54] Daoye Wang, Prashanth Chandran, Gaspard Zoss, Derek Bradley, and Paulo Gotardo. MORf: Morphable radiance fields for multiview neural head modeling. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. [1](#)
- [55] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 10039–10049, 2021. [3](#)
- [56] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9168–9178, 2021. [5](#)
- [57] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *European Conference on Computer Vision Workshops*, 2018. [5](#), [13](#)
- [58] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022. [3](#)
- [59] Olivia Wiles, A Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *European Conference on Computer Vision*, pages 670–686, 2018. [3](#)
- [60] Yue Wu, Yu Deng, Jiaolong Yang, Fangyun Wei, Qifeng Chen, and Xin Tong. AniFaceGAN: Animatable 3d-aware face image generation for video avatars. In *Advances in Neural Information Processing Systems*, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#)
- [61] Jianfeng Xiang, Jiaolong Yang, Yu Deng, and Xin Tong. GRAM-HD: 3d-consistent image generation at high resolution with generative radiance manifolds. In *IEEE/CVF International Conference on Computer Vision*, 2023. [2](#), [3](#), [5](#), [7](#)
- [62] Hongyi Xu, Guoxian Song, Zihang Jiang, Jianfeng Zhang, Yichun Shi, Jing Liu, Wanchun Ma, Jiashi Feng, and Linjie Luo. OmniAvatar: Geometry-guided controllable 3d head synthesis. In *IEEE/CVF International Conference on Computer Vision*, 2023. [1](#), [2](#), [3](#)
- [63] Sicheng Xu, Jiaolong Yang, Dong Chen, Fang Wen, Yu Deng, Yunde Jia, and Tong Xin. Deep 3d portrait from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7710–7720, 2020. [3](#)
- [64] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *IEEE/CVF International Conference on Computer Vision*, pages 9459–9468, 2019. [3](#)
- [65] Jianfeng Zhang, Zihang Jiang, Dingdong Yang, Hongyi Xu, Yichun Shi, Guoxian Song, Zhongcong Xu, Xinchao Wang, and Jiashi Feng. Avatargen: a 3d generative model for animatable human avatars. In *European Conference on Computer Vision Workshops*, pages 668–685, 2022. [2](#), [4](#)
- [66] Xiaoming Zhao, Fangchang Ma, David Güera, Zhile Ren, Alexander G Schwing, and Alex Colburn. Generative multiplane images: Making a 2d gan 3d-aware. In *European Conference on Computer Vision*, pages 18–35, 2022. [2](#)



Figure 9: Uncurated portrait generation results from our method.

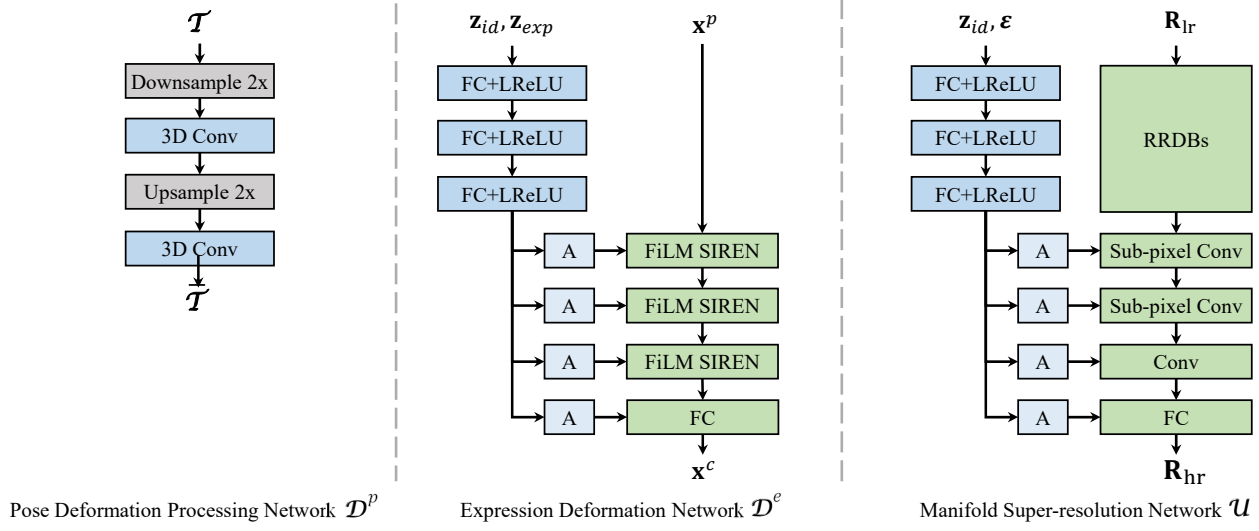


Figure 10: Network structures. \mathcal{D}^p has two 3D conv layers with kernel size $9 \times 9 \times 5$. Both \mathcal{D}^e and \mathcal{U} have a StyleGAN-like structure [25] with a mapping MLP network and a backbone. The backbone of \mathcal{D}^e is mainly based on 3 FiLM SIREN layers [6] while \mathcal{U} uses 4 Residual-in-Residual Dense Blocks (RRDBs) [57] and 2 sub-pixel conv layers [46].