# Generalizable One-shot Neural Head Avatar

**Xueting Li, Shalini De Mello, Sifei Liu, Koki Nagano, Umar Iqbal, Jan Kautz**
NVIDIA

https://research.nvidia.com/labs/lpr/one-shot-avatar

## Abstract

We present a method that reconstructs and animates a 3D head avatar from a single-view portrait image. Existing methods either involve time-consuming optimization for a specific person with multiple images, or they struggle to synthesize intricate appearance details beyond the facial region. To address these limitations, we propose a framework that not only generalizes to unseen identities based on a single-view image without requiring person-specific optimization, but also captures characteristic details within and beyond the face area (*e.g.* hairstyle, accessories, *etc.*). At the core of our method are three branches that produce three tri-planes representing the coarse 3D geometry, detailed appearance of a source image, as well as the expression of a target image. By applying volumetric rendering to the combination of the three tri-planes followed by a super-resolution module, our method yields a high fidelity image of the desired identity, expression and pose. Once trained, our model enables efficient 3D head avatar reconstruction and animation via a single forward pass through a network. Experiments show that the proposed approach generalizes well to unseen validation datasets, surpassing SOTA baseline methods by a large margin on head avatar reconstruction and animation.

## 1 Introduction

Head avatar animation [63, 32, 42] aims to animate a source portrait image with the motion (i.e., pose and expression) from a target image. It is a long-standing task in computer vision that has been widely applied to video conferencing, computer games, Virtual Reality (VR) and Augmented Reality (AR). In real-world applications, synthesizing a realistic portrait image that matches the given identity and motion raises two major challenges – efficiency and high fidelity. Efficiency requires the model to generalize to arbitrary unseen identities and motion without any further optimization during inference. High fidelity demands the model to not only faithfully preserve intricate details in the input image (*e.g.* hairstyle, glasses, earrings), but also hallucinate plausibly whenever necessary (*e.g.* synthesize the occluded facial region when the input is in profile view or generate teeth when the mouth transitions from closed to open).

Traditional methods [16, 19, 13] based on 3D Morphable Models (3DMMs) learn networks that predict shape, expression, pose and texture of an arbitrary source portrait image efficiently. However, these approaches often fall short in synthesizing realistic details due to limited mesh resolution and a coarse texture model. Additionally, they exclusively focus on the facial region while neglecting other personal characteristics such as hairstyle or glasses. Inspired by the remarkable progress made in Generative Adversarial Networks (GANs) [24, 30, 64], another line of methods [63, 72, 75, 54, 51, 78, 17] represent motion as a warping field that transforms the given source image to match the desired pose and expression. Yet, without explicit 3D understanding of the given portrait image, these methods can only rotate the head within limited angles, before exhibiting warping artifacts, unrealistic distortions and undesired identity changes across different target views. Recently, neural rendering [43] has demonstrated impressive results in facial avatar reconstruction and

animation [21, 46, 60, 61, 81, 2, 47, 22, 23, 26, 4]. Compared to meshes with fixed and pre-defined topology, an implicit volumetric representation is capable of learning photo-realistic details including areas beyond the facial region. However, these models have limited capacity and cannot generalize trivially to unseen identities during inference. As a result, they require time-consuming optimization and extensive training data of a specific person to faithfully reconstruct their 3D neural avatars.

In this paper, we present a framework aiming at a more practical but challenging scenario – given an unseen single-view portrait image, we reconstruct an implicit 3D head avatar that not only captures photo-realistic details within and beyond the face region, but also is readily available for animation without requiring further optimization during inference. To this end, we propose a framework with three branches that disentangle and reconstruct the coarse geometry, detailed appearance and expression of a portrait image, respectively. Specifically, given a source portrait image, our *canonical branch* reconstructs its coarse 3D geometry by producing a canonicalized tri-plane [9, 10] with a neutral expression and frontal pose. To capture the fine texture and characteristic details of the input image, we introduce an *appearance branch* that utilizes the depth rendered from the canonical branch to create a second tri-plane by mapping pixel values from the input image onto corresponding positions in the canonicalized 3D space. Finally, we develop an *expression branch* that takes the frontal rendering of a 3DMM with a target expression and a source identity as input. It then produces a third tri-plane that modifies the expression of the reconstruction as desired. After combining all three tri-planes by summation, we carry out volumetric rendering followed by a super-resolution block and produce a high-fidelity facial image with source identity as well as target pose and expression. Our model is learned with large numbers of portrait images of various identity and motion during training. At inference time, it can be readily applied to an unseen single-view image for 3D reconstruction and animation, eliminating the need for additional test-time optimization.

To summarize, our contributions are:

- We propose a framework for 3D head avatar reconstruction and animation that simultaneously captures intricate details in a portrait image while generalizing to unseen identities without test-time optimization.

- To achieve this, we introduce three novel modules for coarse geometry, detailed appearance as well as expression disentanglement and modeling, respectively.

- Our model can be directly applied to animate an unseen image during inference efficiently, achieving favorable performance against state-of-the-art head avatar animation methods.

## 2   Related Works

### 2.1   3D Morphable Models

Reconstructing and animating 3D faces from images has been a fundamental task in computer vision. Following the seminal work by Parke *et al*. [48], numerous methods have been proposed to represent the shape and motion of human faces by 3D Morphable Models (3DMMs) [50, 39, 18, 7, 38]. These methods represent the shape, expression and texture of a given person by linearly combining a set of bases using person-specific parameters. Building upon 3DMMs, many works have been proposed to reconstruct and animate human faces by estimating the person-specific parameters given a single-view portrait image [16, 19, 13, 37]. While 3DMMs provide a strong prior for understanding of human faces, they are inherently limited in two ways. First, they exclusively focus on the facial region and fail to capture other characteristic details such as hairstyle, eye glasses, inner mouth *etc*. Second, the geometry and texture fidelity of the reconstructed 3D faces are limited by mesh resolution, leading to unrealistic appearance in the rendered images. In this work, we present a method that effectively exploits the strong prior in 3DMMs while addressing its geometry and texture fidelity limitation by employing neural radiance fields [43, 5, 45].

### 2.2   2D Expression Transfer

The impressive performance of Generative Adversarial Networks (GANs) [24] spurred another line of head avatar animation methods [63, 72, 75, 54, 51, 78, 17]. Instead of reconstructing the underline 3D shape of human faces, these methods represent motion (*i.e*. expression and pose) as a warping field. Expression transfer is carried out by applying a warping operation onto the source image to

match the motion of the driving image. By leveraging the powerful capacity of generative models, these methods produce high fidelity results with more realistic appearance compared to 3DMM-based methods. However, without an explicit understanding and modeling of the underlying 3D geometry of human faces, these methods usually suffer from warping artifacts, unrealistic distortions and undesired identity change when the target pose and expression are significantly different from the ones in the source image. In contrast, we explicitly reconstruct the underline 3D geometry and texture of a portrait image, enabling our method to produce more realistic synthesis even in cases of large pose change during animation.

## 2.3 Neural Head Avatars

Neural Radiance Field (NeRF) [43, 5, 45] debuts remarkable performance for 3D scene reconstruction. Many works [21, 46, 60, 61, 81, 2, 47, 22, 23, 26, 4, 34, 82, 3] attempt to apply NeRF to human portrait reconstruction and animation by extending it from static scenes to dynamic portrait videos. Although these methods demonstrate realistic reconstruction results, they inefficiently learn separate networks for different identities and require thousands of frames from a specific individual for training. Another line of works focus on generating a controllable 3D head avatar from random noise [59, 67, 57, 44, 68, 36, 84, 56]. Intuitively, 3D face reconstruction and animation could be achieved by combining these generative methods with GAN inversion [52, 20, 70, 66]. However, the individual optimization process for each frame during GAN inversion is computationally infeasible for real-time performance in applications such as video conferencing. Meanwhile, several works [62, 74, 6, 15] focus on reconstructing 3D avatars from arbitrary input images, but they cannot animate or reenact these avatars. Closest to our problem setting, few works explore portrait reconstruction and animation in a few-shot [76] or one-shot [27, 32, 42, 85, 40] manner. Specifically, the ROME method [32] combines a learnable neural texture with explicit FLAME meshes [39] to reconstruct a 3D head avatar, encompassing areas beyond the face region. However, using meshes as the 3D shape representation prevents the model from producing high-fidelity geometry and appearance details. Instead of using explicit meshes as 3D representation, the HeadNeRF [27] and MofaNeRF methods learn implicit neural networks that take 3DMM parameters (*i.e.* identity and expression coefficients or albedo and illumination parameters) as inputs to predict the density and color for each queried 3D point. Additionally, the OTAvatar [42] method proposes to disentangle latent style codes from a pre-trained 3D-aware GAN [9] into separate motion and identity codes, enabling facial animation by exchanging the motion codes. Nonetheless, all three models [27, 85, 42] require laborious test-time optimization, and struggle to reconstruct photo-realistic texture details of the given portrait image presumably because they encode the appearance using a compact latent vector. In this paper, we propose the first 3D head neural avatar animation work that not only generalizes to unseen identities without test-time optimization, but also captures intricate details from the given portrait image, surpassing all previous works in quality.

## 3 Method

We present a framework that takes a source image $I_s$ together with a target image $I_t$ as inputs, and synthesizes an image $I_o$ that combines the identity from the source image and the motion (*i.e.*, expression and head pose) from the target image. The overview of the proposed method is illustrated in Fig. 1. Given a source image including a human portrait, we begin by reconstructing the coarse geometry and fine-grained person-specific details via a canonical branch and an appearance branch, respectively. To align this reconstructed 3D neural avatar with the expression in the target image, we employ an off-the-shelf 3DMM [16] model [1] to produce a frontal-view rendering that combines the identity from the source image with the expression from the target image. Our expression branch then takes this frontal-view rendering as input and outputs a tri-plane that aligns the reconstructed 3D avatar to the target expression. By performing volumetric rendering from the target camera view and applying a super-resolution block, we synthesize a high-fidelity image with the desired identity and motion. In the following, we describe the details of each branch in our model, focusing on answering three questions: a) how to reconstruct the coarse shape and texture of a portrait image with neutral expression in Sec. 3.1; b) how to capture appearance details in the source image in Sec. 3.2; and c) how to model and transfer expression from the target image onto the source image in Sec. 3.3.

---

[1]We introduce the preliminary of the 3DMM in the appendix.

target image $I_t$    source image $I_s$

depth $D_s$    point cloud

2D features $F$    neural point cloud

(b) Lifting

neural point cloud    tri-plane

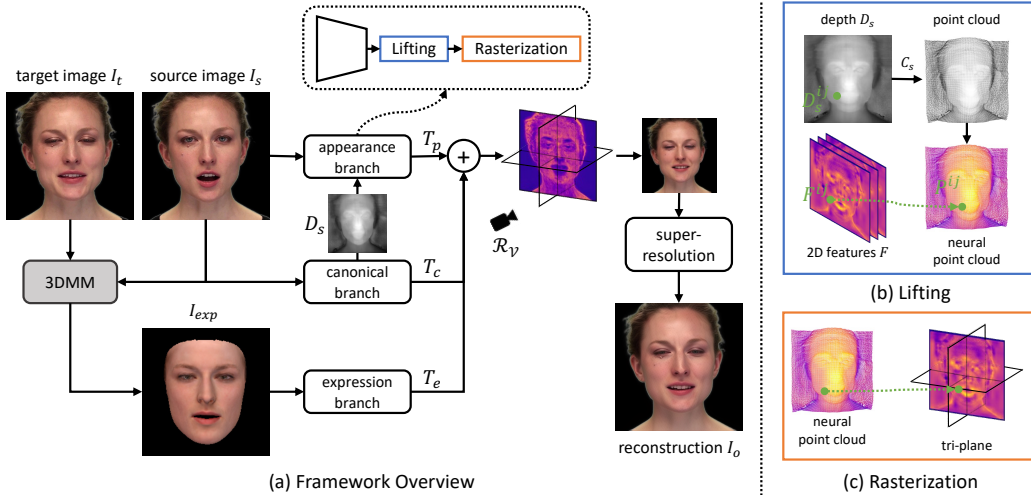(a) Framework Overview

(c) Rasterization

Figure 1: **Overview.** The proposed method contains four main modules: a canonical branch that reconstructs the coarse geometry and texture of a portrait with a neutral expression (Sec. 3.1), an appearance branch that captures fine-grained person-specific details (Sec. 3.2), an expression branch that modifies the reconstruction to desired expression, and a super-resolution block that renders high-fidelity synthesis (Sec. 3.4).

The super-resolution module and the training stages with associated objectives will be discussed in Sec. 3.4 and Sec. 3.5, respectively.

## 3.1   Coarse Reconstruction via the Canonical Branch

Given a source image $I_s$ depicting a human portrait captured from the camera view $C_s$, the canonical branch predicts its coarse 3D reconstruction represented as a tri-plane [9, 10] $T_c$. To serve as a strong geometric prior for the subsequent detailed appearance and expression modeling, we impose two crucial properties on the coarse reconstruction. First, the coarse reconstruction of face images captured from different camera views should be aligned in the 3D canonical space, allowing the model to generalize to single-view portrait images captured from arbitrary camera views. Second, we enforce the coarse reconstruction to have a *neutral* expression (*i.e.*, opened eyes and closed mouth), which facilitates the expression branch to add the target expression effectively.

Based on these two goals, we design an encoder $E_c$ that takes the source image $I_s \in \mathbb{R}^{3 \times 512 \times 512}$ as input and predicts a canonicalized tri-plane $T_c \in \mathbb{R}^{3 \times 32 \times 256 \times 256}$. Specifically, we fine-tune a pre-trained SegFormer [65] model as our encoder, whose transformer design enables effective mapping from the 2D input to the canonicalized 3D space. Furthermore, to ensure that $T_c$ has a neutral expression, we employ a 3DMM [16] to render a face with the same identity and camera pose of the source image, but with a neutral expression. We then encourage the rendering of $T_c$ to be close to the 3DMM's rendering within the facial region by computing an L1 loss and a perceptual loss [28, 77] between them:

$$
\begin{aligned}
I_c &= \mathcal{R}_{\mathcal{V}}(T_c, C_s) \\
I_{neu}, M_{neu} &= \mathcal{R}_{\mathcal{M}}(\alpha_s, \beta_0, C_s) \\
\mathcal{L}_{neutral} &= ||I_{neu} - I_c \times M_{neu}|| + ||\phi(I_{neu}) - \phi(I_c \times M_{neu})||,
\end{aligned}
\tag{1}
$$

where $\mathcal{R}_{\mathcal{V}}(T, C)$ is the volumetric rendering of a tri-plane $T$ from the camera view $C$, $\phi$ is a pre-trained VGG-19 network [55]. $I, M = \mathcal{R}_{\mathcal{M}}(\alpha, \beta, C)$ is the 3DMM [16] that takes identity coefficients $\alpha$, expression coefficients $\beta$ as inputs, and renders an image $I$ and a mask $M$ including only the facial region from camera view $C$. By setting $\alpha = \alpha_s$ (*i.e.* the identity coefficents of $I_s$) and $\beta_0 = \mathbf{0}$ in Eq. 1, we ensure that $I_{neu}$ has the same identity of $I_s$ but with a neutral expression.

As shown in Fig. 2(c), the rendered image $I_c$ from the canonical tri-plane $T_c$ indeed has a neutral expression with opened eyes and closed mouth, but lacks fine-grained appearance. This is because mapping a portrait from the 2D input to the canonicalized 3D space is a challenging and holistic process. As a result, the encoder primarily focuses on aligning inputs from different camera views and neglects individual appearance details. Similar observations have also been noted in [15, 74]. To

4

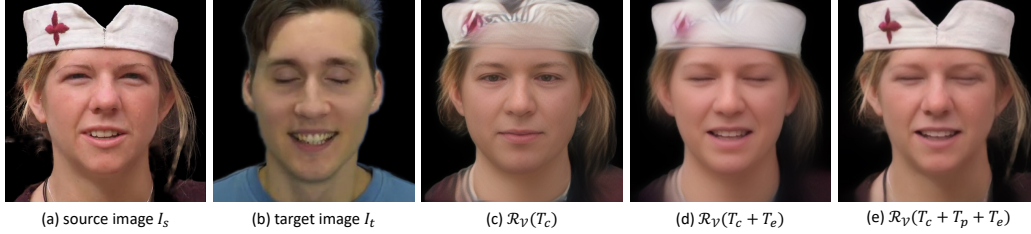| (a) source image $I_s$ | (b) target image $I_t$ | (c) $\mathcal{R}_\mathcal{V}(T_c)$ | (d) $\mathcal{R}_\mathcal{V}(T_c + T_e)$ | (e) $\mathcal{R}_\mathcal{V}(T_c + T_p + T_e)$ |

Figure 2: **Visualization of the contribution of each branch.** (a) Source image. (b) Target image. (c) Rendering of the canonical tri-plane. (d) Rendering of the combination of the canonical and expression tri-planes. (e) Rendering of the combination of all three tri-planes.

resolve this issue, we introduce an appearance branch that spatially transfers details from the input image to the learned coarse reconstruction's surface in the next section.

## 3.2 Detail Reconstruction via the Appearance Branch

We now introduce the appearance branch that aims to capture and reconstruct intricate facial details in the input image. The core idea is to leverage the depth map rendered from the canonical tri-plane $T_c$ to compute the 3D position of each pixel in the image such that the facial details can be accurately "transferred" from the 2D input image to the 3D reconstruction. Specifically, we first render $T_c$ from the source camera view $C_s$ to obtain a depth image $D_s \in \mathbb{R}^{128 \times 128}$. The 3D position (denoted as $P^{ij}$) of each pixel $I_s^{ij}$ in the source image $I_s$ can be computed by $P^{ij} = \mathbf{o} + D_s^{ij}\mathbf{d}$, where $\mathbf{o}$ and $\mathbf{d}$ are the ray origin and viewing direction sampled from the camera view $C_s$ of the source image. Based on the 3D locations of all pixels, we construct a neural point cloud [69, 58] by associating the color information from each pixel $I_s^{ij}$ in the 2D image to its corresponding 3D position $P^{ij}$. Instead of directly using the RGB color of each pixel, we employ an encoder $E_p$ to extract 2D features (denoted as $F \in \mathbb{R}^{32 \times 128 \times 128}$) from $I_s$ and associate the feature at each pixel to its corresponding 3D location. As a result, we establish a neural point cloud composed of all visible pixels in the image and associate each point with a 32-dimensional feature vector. This mapping process from a 2D image to the 3D space, is referred to as "Lifting" and demonstrated in Fig. 1(b).

To integrate the neural point cloud into the canonical tri-plane $T_c$, we propose a "Rasterization" process (see Fig. 1(c)) that converts the neural point cloud to another tri-plane denoted as $T_p$ such that it can be directly added to $T_c$. For each location on the planes (*i.e.* the XY-, YZ-, XZ-plane) in $T_p$, we compute its nearest point in the neural point cloud and transfer the feature from the nearest point onto the query location on the plane. A comparison between Fig. 2(d) and Fig. 2(e) reveals the contribution of our appearance tri-plane $T_p$, which effectively transfers the fine-grained details (*e.g.*, pattern on the hat) from the image onto the 3D reconstruction.

## 3.3 Expression Modeling via the Expression Branch

Expression reconstruction and transfer is a challenging task. Naively predicting the expression from an image poses difficulties in disentangling identity, expression, and head rotation. Meanwhile, 3DMMs provide a well-established expression representation that captures common human expressions effectively. However, the compact expression coefficients in 3DMMs are highly correlated with the expression bases and do not include spatially varying deformation details. As a result, conditioning a network solely on these coefficients for expression modeling can be challenging. Instead, we propose a simple expression branch that fully leverages the expression prior in any 3DMM and seamlessly integrates with the other two branches. The core idea is to provide the model with target expression information using a 2D rendering from the 3DMM instead of the expression coefficients. As shown in Fig. 1(a), given the source image $I_s$ and target image $I_t$, we predict their corresponding shape and expression coefficients denoted as $\alpha_s$ and $\beta_t$ respectively using a 3DMM prediction network [16]. By combining $\alpha_s$ and $\beta_t$, we render a *frontal-view* facial image as $I_{exp} = \mathcal{R}_\mathcal{M}(\alpha_s, \beta_t, C_{front})$, where $C_{front}$ is a pre-defined frontal camera pose. We then use an encoder (denoted as $E_e$) that takes $I_{exp}$ as input and produces an expression tri-plane $T_e \in \mathbb{R}^{3 \times 32 \times 256 \times 256}$. We modify the canonical tri-plane $T_c$ to the target expression by directly adding $T_e$ with $T_c$. Note that we always render $I_{exp}$ in the pre-defined frontal view so that the expression encoder can focus on modeling expression changes only and ignore motion changes caused by head rotation. Moreover, our expression encoder also learns to hallucinate realistic inner mouths (*e.g.*, teeth) according to the target expression, as

the 3DMM rendering $I_{exp}$ does not model the inner mouth region. Fig. 2(d) visualizes the images rendered by combining the canonical and expression tri-planes, where the target expression from Fig. 2(b) is effectively transferred onto Fig. 2(a) through the expression tri-plane.

## 3.4 The Super-resolution Module

By adding the canonical and appearance tri-planes from a source image, together with the expression tri-plane from a target image, we reconstruct and modify the portrait in the source image to match the target expression. Through volumetric rendering, we can obtain a portrait image at a desired camera view. However, the high memory and computational cost of volumetric rendering prevents the model from synthesizing a high-resolution output. To overcome this challenge, existing works [9, 36, 56, 57] utilize a super-resolution module that takes a low-resolution rendered image or feature map as input and synthesizes a high-resolution result. In this work, we follow this line of works and fine-tune a pre-trained GFPGAN [64, 57] as our super-resolution module [57]. By pre-training on the task of 2D face restoration, GFPGAN learns a strong prior for high-fidelity facial image super-resolution. Additionally, its layer-wise feature-conditioning design prevents the model from deviating from the low-resolution input, thereby mitigating temporal or multi-view inconsistencies, as observed in [57].

## 3.5 Model Training

We utilize a two-stage training schedule to promote multi-view consistent reconstructions, as well as to reduce the overall training time. In the first stage, we train our model without the super-resolution module using a reconstruction objective and the neutral expression loss discussed in Sec. 3.1. Specifically, we compare the rendering of a) the canonical tri-plane (*i.e.*, $I_c = \mathcal{R}_\mathcal{V}(T_c, C_t)$), b) the combination of the canonical and expression tri-planes (*i.e.*, $I_{c+e} = \mathcal{R}_\mathcal{V}(T_c + T_e, C_t)$), and c) the combination of all three tri-planes (*i.e.*, $I_{c+e+p} = \mathcal{R}_\mathcal{V}(T_c + T_e + T_p, C_t)$) with the target image via the L1 and the perceptual losses similarly to Eq. 1:

$$\begin{aligned} \mathcal{L}_1 &= ||I_c - I_t|| + ||I_{c+e} - I_t|| + ||I_{c+e+p} - I_t||, \\ \mathcal{L}_p &= ||\phi(I_c) - \phi(I_t)|| + ||\phi(I_{c+e}) - \phi(I_t)|| + ||\phi(I_{c+e+p}), \phi(I_t)||. \end{aligned} \quad (2)$$

Intuitively, applying supervision to different tri-plane combinations encourages the model to predict meaningful reconstruction in all three branches. To encourage smooth tri-plane reconstruction, we also adopt the TV loss proposed in [9]. The training objective for the first stage is $\mathcal{L}^1 = \lambda_1 \mathcal{L}_1 + \lambda_p \mathcal{L}_p + \lambda_{TV} \mathcal{L}_{TV} + \lambda_{neutral} \mathcal{L}_{neutral}$, where $\lambda_x$ is the weight of the corresponding objective.

In the second stage, to encourage multi-view consistency, we only fine-tune the super-resolution module and freeze other parts of the model. We use all the losses in the first stage and a dual-discriminator proposed in [9] that takes the concatenation of the low-resolution rendering and the high-resolution reconstruction as input. Specifically, we use the logistic formulation [24, 30, 64] of the adversarial loss $L_{adv} = \mathbb{E}_{(I_o^l, I_o)} \text{softplus}(D(I_o^l \oplus I_o)))$, where $I_o^l$ is the upsampled version of the low-resolution rendered image and $\oplus$ represents the concatenation operation. The overall training objective of the second stage is $\mathcal{L}^2 = \mathcal{L}^1 + \lambda_{adv} \mathcal{L}_{adv}$.

# 4 Experiments

## 4.1 Datasets

**Training datasets.** We train our model using a single-view image dataset (FFHQ [29]) and two video datasets (CelebV-HQ [83] and RAVDESS [41]). For the single-view images in FFHQ, we carry out the 3D portrait reconstruction task, *i.e.*, the source and target images are exactly the same. For the CelebV-HQ and the RAVDESS datasets, we randomly sample two frames from the same video to formulate a pair of images with the same identity but different motion. Furthermore, we observe that parts of videos in the CelebV-HQ and the RAVDESS datasets are fairly static, leading to a pair of source and target images with similar head poses, impeding the model from learning correct 3D shape of portraits. To enhance the learning of 3D reconstruction, we further employ an off-the-shelf 3D-aware GAN model [9] to synthesize 55,857 pairs of images rendered from two randomly sampled camera views, *i.e.*, the source and target images have the same identity and expression but different views.
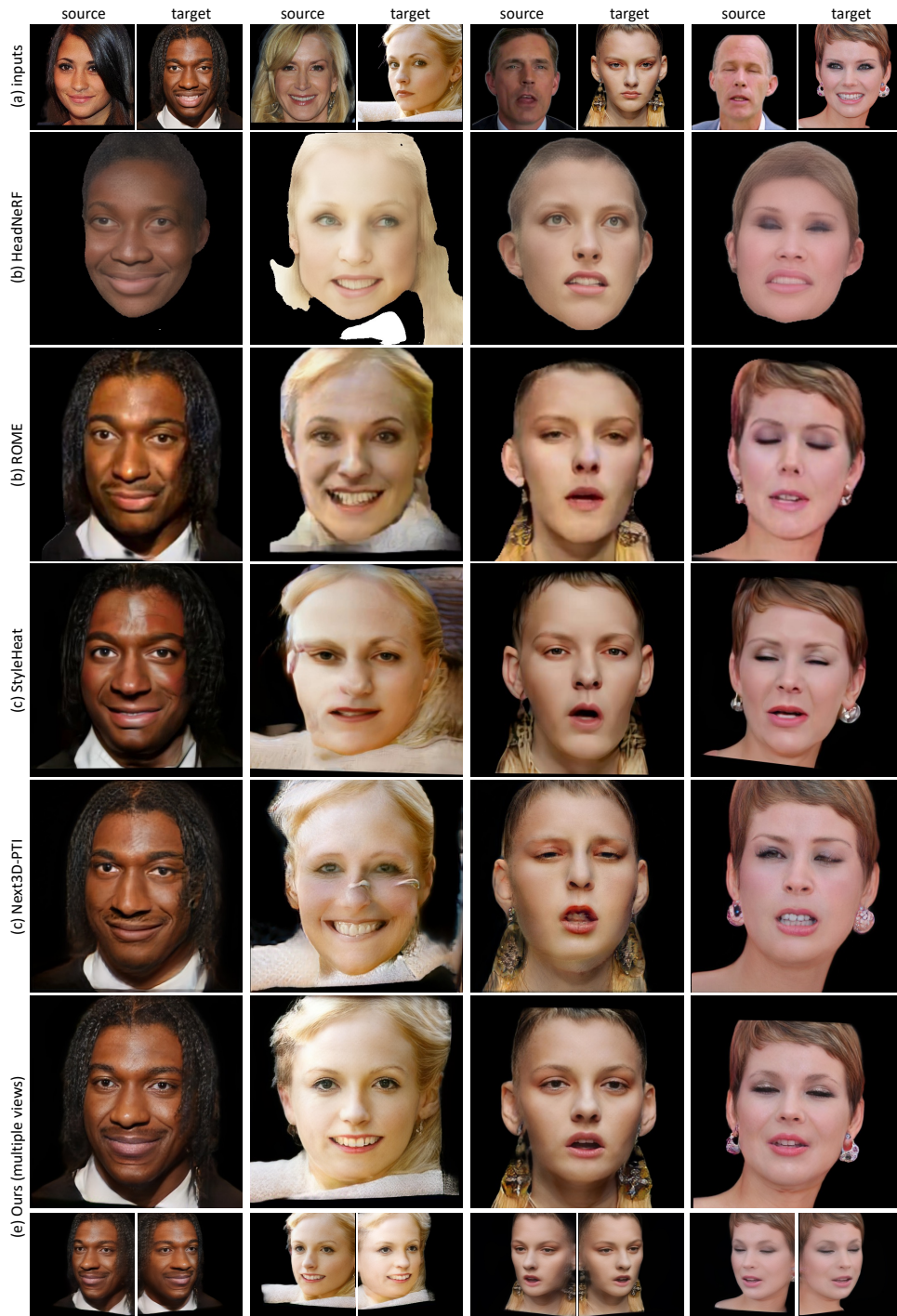
Figure 3: **Cross-identity reenactment on CelebA [35] and HDTF [79].** The first two rows show cross-identity reenactment results on the CelebA dataset, while the last two rows demonstrate motion transfer from videos in the HDTF dataset to images in the CelebA dataset.

**Evaluation datasets.** We evaluate our method and the baselines [32, 27, 42, 72] on the CelebA dataset [35] and the testing split of the HDTF dataset [79], following [72, 42]. Note that our method has never seen any image from these two datasets during training while both StyleHeat[72] and OTAvatar [42] are trained using the training split of the HDTF dataset. Nonetheless, our method generalizes well to all validation datasets and achieves competitive performance, as discussed later.

7

Table 1: **Comparison on CelebA [35].** [†] Evaluated on a subset of CelebA, as discussed in Sec. 4.4.

| Methods | 3D Portrait Reconstruction | | | | | Cross-Identity Reeanct | | | |
|---|---|---|---|---|---|---|---|---|---|
| | L1↓ | LPIPS↓ | PSNR↑ | SSIM↑ | FID↓ | CSIM↑ | AED↓ | APD↓ | FID↓ |
| ROME [32] | 0.032 | 0.085 | 23.47 | 0.847 | 11.00 | 0.505 | **0.244** | 0.032 | 34.45 |
| Ours | **0.015** | **0.040** | **28.61** | **0.946** | **2.457** | **0.531** | 0.251 | **0.023** | **25.26** |
| HeadNeRF[†] [27] | 0.135 | 0.314 | 13.86 | 0.748 | 65.87 | 0.224 | 0.285 | 0.027 | 117.1 |
| Ours[†] | **0.024** | **0.098** | **25.83** | **0.883** | **9.400** | **0.591** | **0.278** | **0.017** | **22.97** |

**Datasets pre-processing.** We compute the camera poses of images in all training and testing datasets using [16] following [9]. As in previous literature [32, 27, 85], background modeling is out of the scope of this work; we further use an off-the-shelf portrait matting method [31] to remove backgrounds in all training and testing images.

## 4.2 Metrics and Baselines

We evaluate all methods for 3D portrait reconstruction, same-identity and cross-identity reenactment.

**3D portrait reconstruction.** We use all 29,954 high-fidelity images[2] in CelebA [35]. We measure the performance by computing various metrics between the reconstructed images and the input images, including the L1 distance, perceptual similarity metric (LPIPS), peak signal-to-noise ratio (PSNR), structural similarity index (SSIM), and Fréchet inception distance (FID).

**Same-identity reenactment.** We follow [42] and use the testing split of HDTF [79], which includes 37,860 frames in total. We use the first frame of each video as the source image and the rest of the frames as target images. Following the evaluation protocol in [42, 72], we evaluate PSNR, SSIM, cosine similarity of the identity embedding (CSIM) based on [14], average expression distance (AED) and average pose distance (APD) based on [16], average keypoint distance (AKD) based on [8], as well as LPIPS, L1 and FID between the reenacted and ground truth frames.

**Cross-identity reenactment.** We conduct cross-identity reenactment on both the CelebA and HDTF datasets. For CelebA, we split the dataset into 14,977 image pairs and transfer the expression and head pose from one image to the other. As for the HDTF dataset, we follow [42] and use one clip as the driving video and the first frame of the other videos as source images, which produces 67,203 synthesized images in total. To fully evaluate the performance of *one-shot* avatar animation, we transfer motion from the HDTF videos to the single-view images in the CelebA dataset. Similar to [72], we use the first 100 frames in the HDTF videos as target images and 60 images sampled from the CelebA as source images, resulting in a total of 114,000 synthesized images. Since there is no ground truth for cross-identity reenactment, we evaluate the results based on the CSIM, AED, APD, and FID metrics. More evaluation details can be found in the appendix.

**Baselines.** In terms of baselines, we compare our method against a 2D talking head synthesis method [72], three 3D head avatar animation methods [27, 32, 42], and a baseline that combines 3D-aware generative model [56] with pivotal tuning [52] (dubbed as "Next3D-PTI" in the following) for 3D avatar reconstruction and animation.

## 4.3 Implementation Details

We implement the proposed method using the PyTorch framework [49] and train it with 8 32GB V100 GPUs. The first training stage takes 6 days, consisting of 750000 iterations. The second training stage takes 2 days with 75000 iterations. Both stages use a batch size of 8 with the Adam optimizer [33] and a learning rate of 0.0001. More implementation details can be found in the appendix.

Table 2: **Comparison on the HDTF dataset [79].**

| Methods | Same-Identity Reenactment | | | | | | | | | Cross-Identity Reenactment | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | CSIM↑ | AED↓ | APD↓ | AKD↓ | LPIPS↓ | L1↓ | FID↓ | CSIM↑ | AED↓ | APD↓ | FID↓ |
| ROME [32] | 20.75 | 0.838 | 0.746 | **0.123** | 0.012 | 2.938 | 0.173 | 0.047 | 31.55 | 0.629 | 0.247 | 0.020 | **43.38** |
| OTAvatar [42] | 20.12 | 0.806 | 0.619 | 0.162 | 0.017 | 2.933 | 0.198 | 0.053 | 36.63 | 0.514 | 0.282 | 0.028 | 44.86 |
| StyleHeat [72] | 19.18 | 0.805 | 0.654 | 0.141 | 0.021 | 2.843 | 0.194 | 0.056 | 108.3 | 0.537 | **0.246** | 0.025 | 105.1 |
| Next3D-PTI [56] | 19.89 | 0.813 | 0.645 | 0.137 | 0.035 | **1.449** | 0.180 | 0.053 | 41.66 | 0.581 | 0.291 | 0.045 | 101.8 |
| Ours | **22.15** | **0.868** | **0.789** | 0.129 | **0.010** | 2.596 | **0.117** | **0.037** | **21.60** | **0.643** | 0.263 | **0.018** | 47.39 |

## 4.4 Qualitative and Quantitative Results

**3D portrait reconstruction.** For each testing portrait image, we reconstruct its 3D head avatar and render it from the observed view using different methods. By comparing the rendering with the input image, we assess the fidelity of 3D reconstruction of each method. Table 1 shows the quantitative results, demonstrating that our model achieves significantly better reconstruction and fidelity scores. These results highlight the ability of our model to faithfully capture details in the input images and reconstruct high-fidelity 3D head avatars. Visual examples are present in the appendix.

**Cross-identity reenactment.** Fig. 3, and Fig. 4 showcase the qualitative results of cross-identity reenactment on the CelebA [35] and HDTF [79] dataset. Compared to the baselines [32, 72, 42], our method faithfully reconstructs intricate details such as hairstyles, earrings, eye glasses *etc*. in the input portrait images. Moreover, our method successfully synthesizes realistic appearance change corre-

Table 3: Cross-identity reenactment between the HDTF dataset [79] and the CelebA dataset [35].

| Methods | CSIM↑ | AED↓ | APD↓ | FID↓ |
|---|---|---|---|---|
| ROME [32] | 0.521 | 0.270 | 0.022 | 76.03 |
| StyleHeat [72] | 0.461 | 0.270 | 0.038 | 94.28 |
| Next3D-PTI [56] | 0.483 | **0.266** | 0.042 | **56.01** |
| Ours | **0.551** | 0.274 | **0.017** | 59.48 |

sponding to the target motion. For instance, our model is able to synthesize plausible teeth when the mouth transitions from closed to open (*e.g.*, row 3 in Fig. 3), it also hallucinates the occluded face region when the input image is in profile view (*e.g.*, row 2 in Fig. 3). In contrast, the mesh-based baseline [32] can neither capture photo-realistic details nor hallucinate plausible inner mouths, while the 2D talking head synthesis baseline [72] produces unrealistic warping artifacts when the input portrait is in the profile view (*e.g.* row 2 in Fig. 3). We provide quantitative evaluations of the cross-identity reenactment results in Table 1, Table 2, and Table 3. Our method demonstrates better fidelity and identity preservation scores, showing its strong ability in realistic portrait synthesis. It is worth noting that HDTF [79] includes images that are less sharp compared to the high-fidelity images our model is trained on, which may account for the slightly lower FID score in Table 2.

**Same-identity reenactment.** Table 2 shows the quantitative results of same-identity reenactment on HDTF [79]. Our method generalizes well to HDTF and achieves better metrics compared to existing SOTA methods [32, 42, 72]. The qualitative results of same-identity reenactment can be found in the appendix.
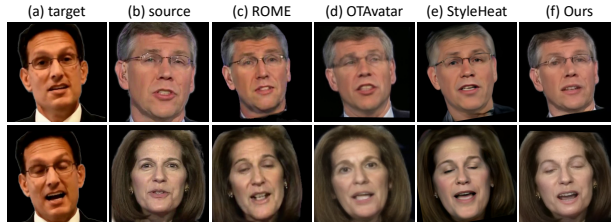


(a) target  (b) source  (c) ROME  (d) OTAvatar  (e) StyleHeat  (f) Ours

Figure 4: **Cross-identity reenactment on HDTF [79].**

**Efficiency.** Since HeadNeRF [27] OTAvatar [42], Next3D-PTI [56] require latent code optimization for unseen identities, the process of reconstructing and animating an avatar takes them 53.0s, 19.4s, 300s, respectively. Meanwhile, ROME [32] and our method only needs an efficient forward pass of the network for unseen identities, taking 1.2s and 0.6s respectively. Overall, our method strikes the best balance in terms of speed and quality.

## 4.5 Ablation Studies

We conduct experiments to validate the effectiveness of the neutral expression constraint (Sec. 3.1), the contribution of the appearance branch (Sec. 3.2), and the design of the expression branch (Sec. 3.3).

---

[2]Due to the time-consuming nature of HeadNeRF, we compare our method with HeadNeRF on a subset of 3000 images from the CelebA dataset.
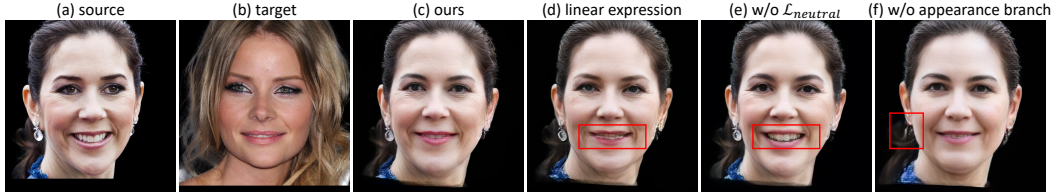
Figure 5: **Ablation studies.** Details are explained in Sec. 4.5.

Table 4: **Ablation studies.** Blue text highlights the inferior performance of the variants. (Sec. 4.5)

| Methods | 3D Portrait Reconstruction | | | | | Cross-Identity Reeanct | | | |
|---|---|---|---|---|---|---|---|---|---|
| | L1↓ | LPIPS↓ | PSNR↑ | SSIM↑ | FID↓ | CSIM↑ | AED↓ | APD↓ | FID↓ |
| w/o appearance | 0.061 | 0.239 | 19.53 | 0.712 | 47.84 | 0.370 | 0.243 | 0.015 | 43.33 |
| w/o neutral constraint | 0.027 | 0.124 | 25.65 | 0.854 | 13.56 | 0.593 | 0.315 | 0.018 | 22.92 |
| linear expression | 0.031 | 0.134 | 25.08 | 0.841 | 14.62 | 0.443 | 0.217 | 0.016 | 26.18 |
| Ours | 0.030 | 0.116 | 24.77 | 0.861 | 10.47 | 0.599 | 0.276 | 0.017 | 17.36 |

**Neutral expression constraint.** In our model, we aim to wipe out the expression in the source image by enforcing the coarse reconstruction from the canonical branch to have a neutral expression. This ensures that the expression branch always animates a fully "neutralized" expressionless face from the canonical branch. Without this, the expression branch fails to correctly modify the coarse reconstruction into the target expression, as shown in Fig. 5(e) and Table 4 (*i.e.*, worse AED score).

**Appearance branch.** The appearance branch is the key to reconstructing intricate facial details of the input portrait image. Without this branch, the model struggles to capture photo-realistic details, resulting in considerably lower reconstruction and fidelity metrics, as shown in Table 4 and Fig. 5(f).

**Alternative expression branch design.** Instead of using the frontal view rendering from the 3DMM to provide target expression information to the expression branch (see Sec. 3.3), an alternative way is to use the 3DMM expression coefficients to linearly combine a set of learnable expression bases. However, this design performs sub-optimally as it overlooks the individual local deformation details caused by expression changes, introducing artifacts (*e.g.*, mouth not fully closed) shown in Fig. 5(d) and lower FID in Table 4. We provide more results and ablations in the appendix.

## 5  Conclusions and Broader Impact

**Conclusions.** In this paper, we propose a framework for one-shot 3D human avatar reconstruction and animation from a single-view image. Our method excels at capturing photo-realistic details in the input portrait image, while simultaneously generalizing to unseen images without the need for test-time optimization. Through comprehensive experiments and evaluations on validation datasets, we demonstrate that the proposed approach achieves favorable performance against state-of-the-art baselines on head avatar reconstruction and animation.

**Broader Impact.** The proposed framework has the potential to make significant contributions to various fields such as video conferencing, entertainment industries, and virtual reality. It offers a wide range of applications, including but not limited to animating portraits for film/game production, or reducing transmission costs in video conferences through self-reenactment that only requires transmitting a portrait image with compact motion vectors. However, the proposed method may present important ethical considerations. One concerning aspect is the possibility of generating "deepfakes", where manipulated video footage portrays individuals saying things they have never actually said. This misuse could lead to serious privacy infringements and the spread of misinformation. We do not advocate for such activities and instead underscore the need to build guardrails to ensure safe use of talking-head technology, such as [53, 80, 12, 1, 11].

# Appendix

In this document, we first show more ablation studies in Sec. A. We further demonstrate more qualitative results in Sec. B. Additional details of the proposed framework and the evaluation process are described in Sec. C and Sec. D. Finally, we discuss preliminaries of the 3DMMs in Sec. E and the limitations of the proposed method in Sec. F, respectively.
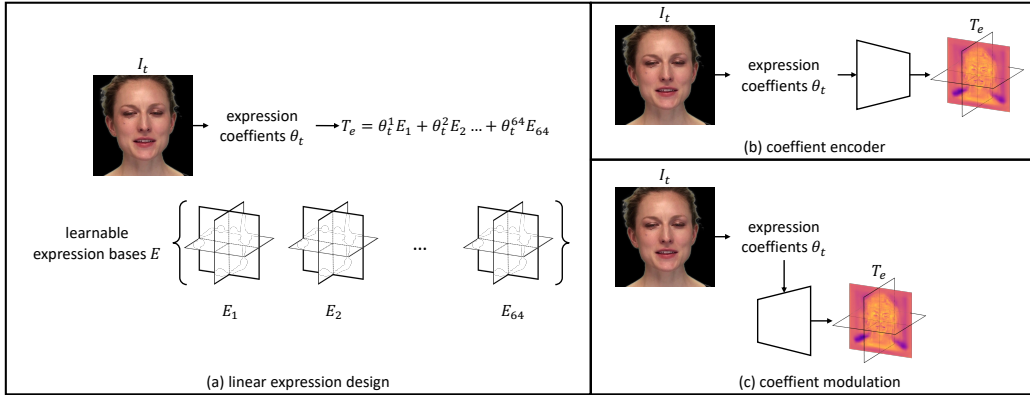
## A   More Ablation Studies



Figure 6: **Ablation variants.** See Sec. A for details.

Table 5: **Ablation studies.** Blue text highlights the inferior performance of the variants.

| Methods | 3D Portrait Reconstruction | | | | | Cross-Identity Reeanct | | | |
|---|---|---|---|---|---|---|---|---|---|
| | L1↓ | LPIPS↓ | PSNR↑ | SSIM↑ | FID↓ | CSIM↑ | AED↓ | APD↓ | FID↓ |
| modulation | 0.030 | 0.130 | 25.01 | 0.841 | 11.80 | 0.558 | 0.280 | 0.018 | 22.84 |
| encoder | 0.028 | 0.121 | 25.69 | 0.850 | 9.840 | 0.538 | 0.264 | 0.017 | 18.88 |
| fine-tune high-res | 0.028 | 0.114 | 25.87 | 0.870 | 9.177 | 0.597 | 0.278 | 0.017 | 21.51 |
| ours | 0.030 | 0.116 | 24.77 | 0.861 | 10.47 | 0.599 | 0.276 | 0.017 | 17.36 |

**Details of the linear expression branch.**    We show the architecture of an alternative expression branch design in Fig. 6 (a). As described in Sec.4.5 of the main submission, this design draws inspirations from 3DMMs and learns 64 tri-plane-based expression bases denoted as $\{E_1, ..., E_{64}\}$. To produce the target expression tri-plane, it linearly combines the learnable bases by $T_c = \theta_t^1 E_1 + ... + \theta_t^{64} E_{64}$, where $\theta_t$ are the target expression coefficients extracted from the target image by the 3DMM [16]. As shown in Fig.5 and Table.4 in the main submission, this design produces unrealistic mouth regions during the animation.

**Expression branch using coefficients.**    An intuitive design for the expression branch is to utilize an encoder that directly maps the target 3DMM expression coefficients to an expression tri-plane. We investigate this design by using two kinds of encoders: i) We simply use linear layers followed by transpose convolutional layers to map the expression coefficient vector to an expression tri-plane. We dub this design as the "encoder" design and show its architecture in Fig. 6(b). ii) We use the generator from StyleGAN2 [30] as our encoder. Specifically, we first map the target expression coefficients to a set of style vectors, the encoder takes a constant tensor as input and modulates the features at each layer using the style vectors produced from the expression coefficients. We denote this design as "modulation" and show its structure in Fig. 6(c).

As shown in Table 5, for the cross-identity reenactment evaluation on CelebA [35], both alternative expression branch designs discussed above have lower CSIM score. This is because the expression coefficients are orthogonal to the identity in the source image. By taking the expression coefficients alone as input, the model has less information of the source identity and suffers from identity preservation while animating.

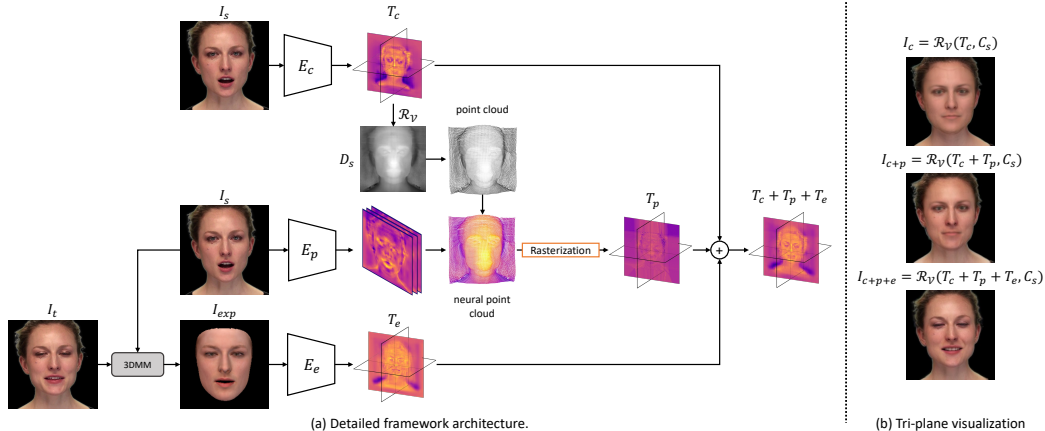(a) Detailed framework architecture.

(b) Tri-plane visualization

Figure 7: **Framework architecture.** (a) We present an extended version of Fig.1(a) in the submission with more framework details. The volumetric rendering and super-resolution block are left out. (b) Visualization of tri-plane combinations.

**Fine-tuning end-to-end in stage II.** As discussed in Sec.3.5 in the submission, in order to preserve multi-view consistency, we only fine-tune the super-resolution module while freezing other parts in Stage II. We verify the effectiveness of this choice by conducting an ablation study where we instead fine-tune end-to-end in Stage II.

Table 5 demonstrates the quantitative evaluation of this variant model. Though it has better reconstruction results from the observed view, it has worse FID score for the task of cross-identity reenactment. This indicates this variant synthesizes less realistic animation results at the target pose.

# B More Qualitative Results

In this section, we show more qualitative results on the CelebA [35] and HDTF [79] datasets.

## B.1 Cross-identity Reenactment on CelebA

In Fig. 8, Fig. 9, Fig. 10 and Fig. 11, we present more visualizations of cross-identity reenactment on the CelebA dataset.

## B.2 3D Portrait Reconstruction on CelebA.

In Fig. 12, Fig. 13, Fig. 14 and Fig. 15, we show portrait reconstruction results by the proposed method visualized in different views .

## B.3 Cross-identity Reenactment from HDTF to CelebA

We visualize more cross-identity reenacment results from the videos in HDTF to the images in CelebA, including comparison with the baselines [32, 72, 56] in Fig. 16, Fig. 17, Fig. 18, Fig. 19 and the supplementary video.

## B.4 Same-identity Reenactment on HDTF

Fig. 20 shows the qualitative results of same-identity reenactment on the HDTF dataset, as well as comparison with OTAvatar [42], ROME [32], StyleHeat [72] and Next3D-PTI [56].

## B.5 Cross-identity Reenactment on HDTF

We present more qualitative results of cross-identity reenactment on the HDTF dataset, as well as comparison with OTAvatar [42], ROME [32], StyleHeat [72] and Next3D-PTI [56] in Fig. 21.

## C  More Implementation Details

**Detailed network architecture.**  In Fig. 7, we present the detailed architecture of the proposed method. As discussed in Sec.3 of the main submission, our model includes three branches that capture the coarse geometry, detailed appearance and expression, respectively. Specifically, the encoder $E_c$ in the canonical branch takes a source image of size $3 \times 512 \times 512$ as input and outputs a feature map of size $256 \times 128 \times 128$. By passing the feature map through four convolution layers and one transpose convolution layer, we obtain a canonical tri-plane of size $3 \times 32 \times 256 \times 256$. The encoder $E_p$ in the appearance branch takes the source image as input and outputs a feature map of size $256 \times 128 \times 128$. Through the "Lifting" and "Raterization" process introduced in Sec.3.2 of the main submission, we produce an appearance tri-plane of size $3 \times 32 \times 256 \times 256$. Furthermore, to prevent the expression in the source image from leaking into the final animation, we use an off-the-shelf face parsing network [73][3] to mask out the eye and mouth regions before providing the source image to the encoder $E_p$. Finally, the encoder $E_e$ in the expression branch is similarly designed as $E_c$, except that it takes the frontal-view 3DMM rendering with the target expression as input. All three encoders (*i.e.* $E_c, E_p, E_e$) use a pre-trained SegFormer [65] model, up to the classifier layer. We adopt the tri-plane decoder proposed by [9] to map the interpolated tri-plane feature to color and density for each 3D point. For the super-resolution block, we fine-tune a pre-trained GFPGAN [64] model without modifying its architecture.

## D  More Evaluation Details

We explain details of how we evaluated the baseline methods and the proposed method in this section.

**3D portrait reconstruction.**  For the ROME method [32], we use the publicly available code and model[4], which renders $256^2$ images. For fair comparison, we resize our prediction and the ground truth images from $512^2$ to $256^2$. Since the synthesis from ROME is not pixel-to-pixel aligned to the input image, we rigidly transform the ground truth image and our image such that they align with ROME's predictions. To this end, we apply the Procrustes process [25, 71] to align our prediction/ground truth image to ROME's prediction using facial landmarks detected by [8]. We also replace the white background in ROME's implementation to a black one to match the ground truth images and our predictions. We compare with ROME on all 29,954 high-fidelity images in CelebA [35] for 3D portrait reconstruction.

Since the synthesis results by HeadNeRF [27][5] and our method are pixel-to-pixel aligned to the input image, we do not carry out further alignment when comparing to HeadNeRF. However, as discussed in Sec.4.2 of the main submission, applying HeadNeRF on all 29,954 images in CelebA is computationally infeasible. Thus we apply our method and HeadNeRF on a randomly sampled subset that includes 3000 images from the CelebA dataset.

**Reenactment.**  We compare with ROME [32], StyleHeat [72][6], OTAvatar [42] and Next3D [56] combined with PTI [52] for same-identity and cross-identity reenactment. We leave HeadNeRF out on the HDTF dataset since it is impossible to test it on tens of thousands frames due to the time-consuming optimization for each frame. Moreover, OTAvatar [42] is a concurrent work and up to the date of this paper, only its partial code[7] that allows for comparison on the HDTF dataset alone has been released publicly. So we do not compare with it on motion transfer from the HDTF dataset to the CelebA dataset. To animate a given portrait, the Next3D method [56] first uses pivotal tuning [52] to map the portrait image to the latent space of its generator and then animates the portrait using expressions extracted from the target video. We use the publicly available implementation[8] of Next3D with pivotal tuning. For fair comparison, we align predictions from all methods to the target image using the Procrustes process discussed above. Note that synthesis from all methods have a black background and are readily comparable after the alignment.

---

[3]https://github.com/zllrunning/face-parsing.PyTorch
[4]https://github.com/SamsungLabs/rome
[5]https://github.com/CrisHY1995/headnerf
[6]https://github.com/FeiiYin/StyleHEAT
[7]https://github.com/theEricMa/OTAvatar
[8]https://github.com/MrTornado24/Next3D

# E    Preliminaries of 3DMMs

We exploit the geometry prior from a 3DMM [50] that represents the shape and texture of a portrait by:

$$S = \bar{S} + B_{id}\alpha + B_{exp}\beta$$
$$T = \bar{T} + B_{tex}\delta \tag{3}$$

where $\bar{S}, \bar{T}$ are the mean shape and texture of human faces, $B_{id}, B_{exp}, B_{tex}$ are the shape, expression and texture bases, and $\alpha, \beta, \delta$ are coefficients that linearly combine the shape, expression and texture bases, respectively. Since we mainly utilize the shape and expression components in the 3DMM in this work, we ignore its texture and illumination modules and simply denote the rendering operation from a camera view $C$ as $I, M = \mathcal{R}_M(\alpha, \beta, C)$, where $I$ is the rendered image, and $M$ is the rendered mask that only includes the facial region.

# F    Limitations

**Teeth and pupil reconstruction.**    3D head avatar reconstruction and animation is a highly challenging task. The proposed method takes the first step to produce high-fidelity results. However, to generalize to any portrait image, one dilemma is that the expression of the source portrait and target image could be arbitrary, which introduces various challenging scenarios. For instance, the source portrait image could have a closed mouth while the target expression has an open mouth (*e.g.* the second row of Fig. 8). In this case, the model should hallucinate correct inner mouth regions. Yet, in other cases, the inner mouth is visible in the source portrait (*e.g.* the sixth row in Fig. 8). To resolve this dilemma, in this work, our model simply always hallucinates the inner mouth of the individual through our expression branch. As a result, the hallucinated mouth could deviate from the one in the source image. In other words, our model cannot accurately reconstruct teeth and lips, as shown in Fig. 22. The same analysis applies to pupils. Since we have no prior knowledge of whether the eyes in the portrait image are open or closed, our model always hallucinates the pupils through the expression branch instead of reconstructing the ones in the source image. We leave this limitation to future works.

source    target   target view     other views

Figure 8: **Cross-identity reenactment on CelebA.**

| source | target | target view | other views |
|--------|--------|-------------|-------------|

Figure 9: **Cross-identity reenactment on CelebA.**
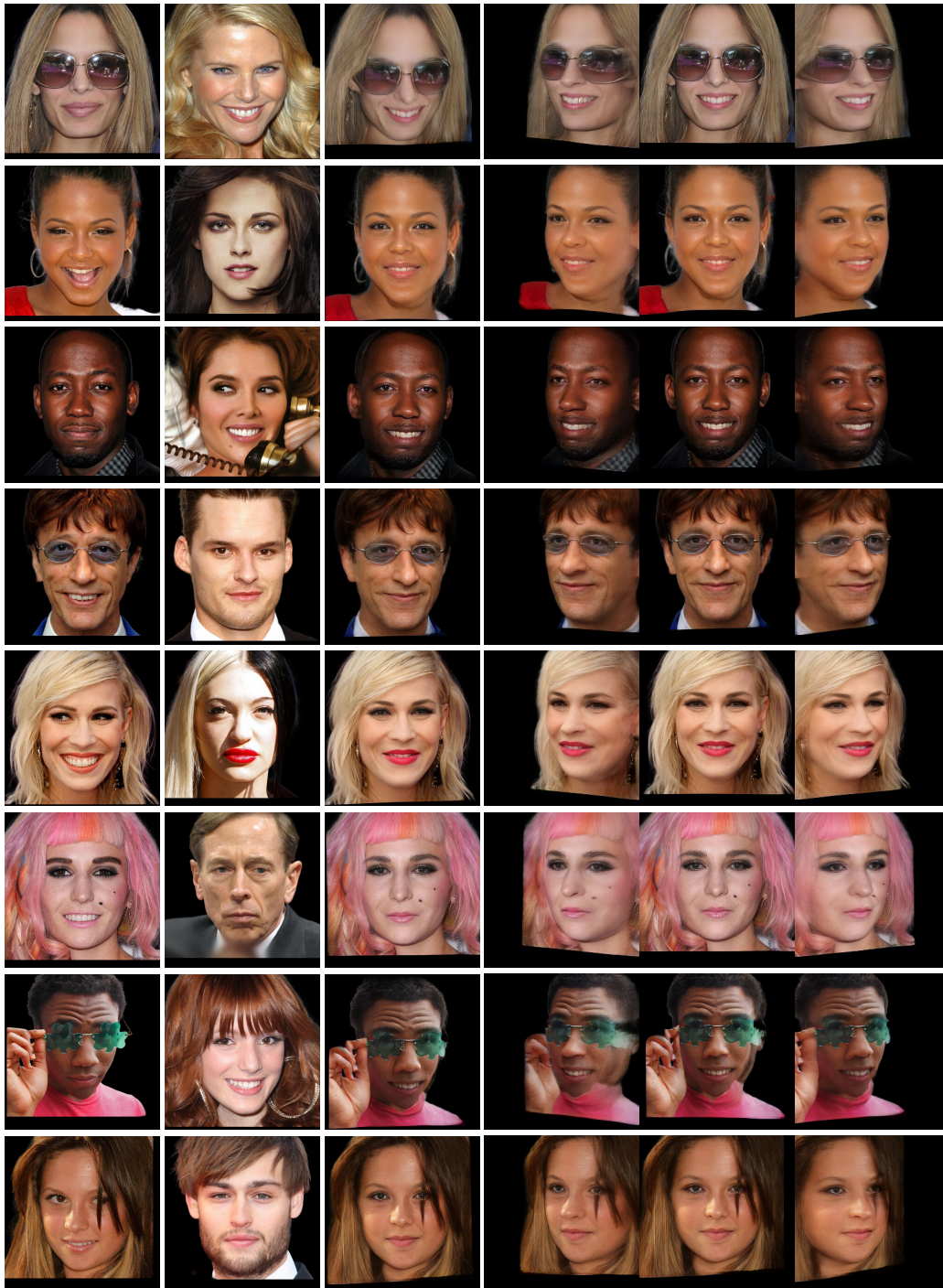
source        target        target view        other views

Figure 10: **Cross-identity reenactment on CelebA.**

source         target       target view        other views
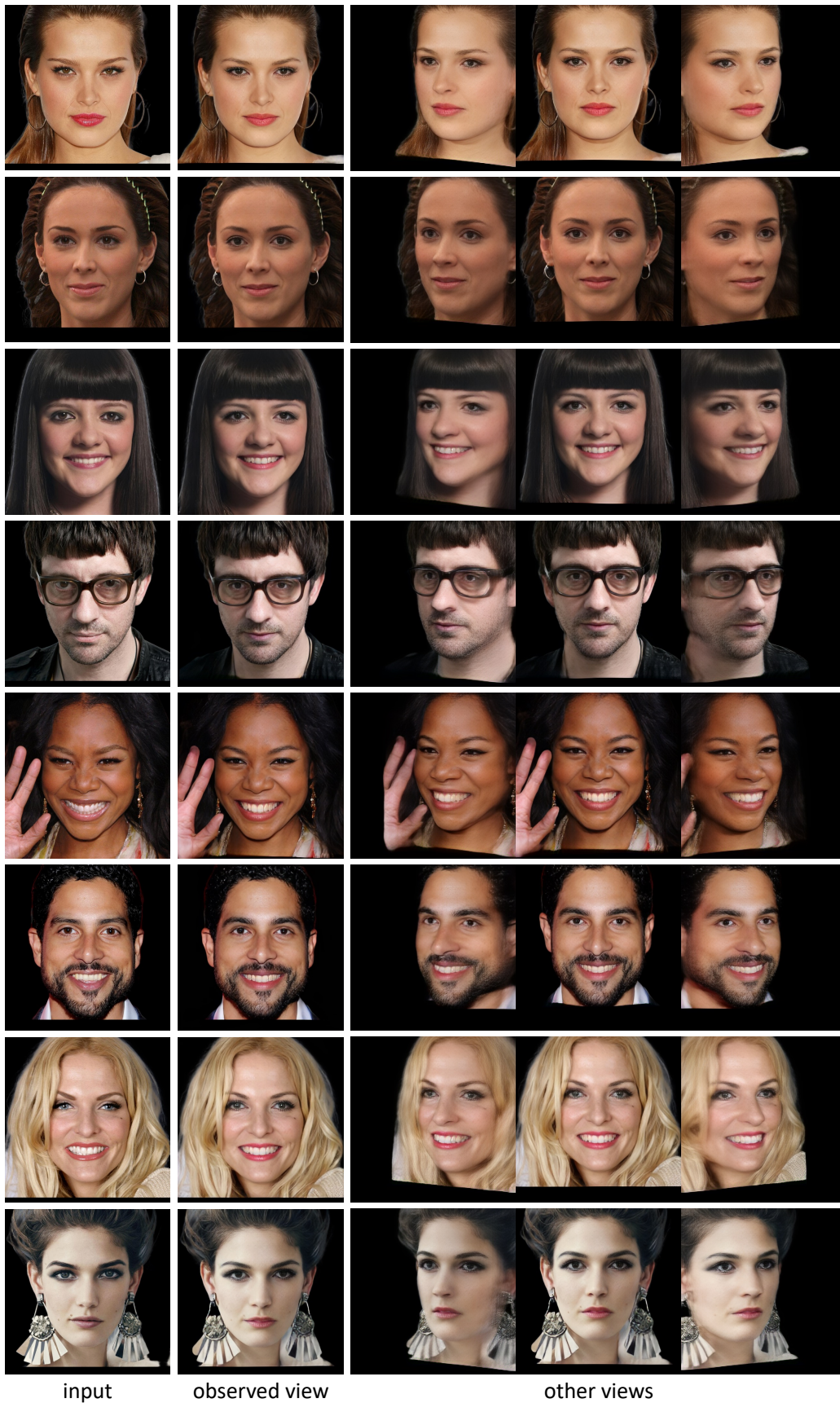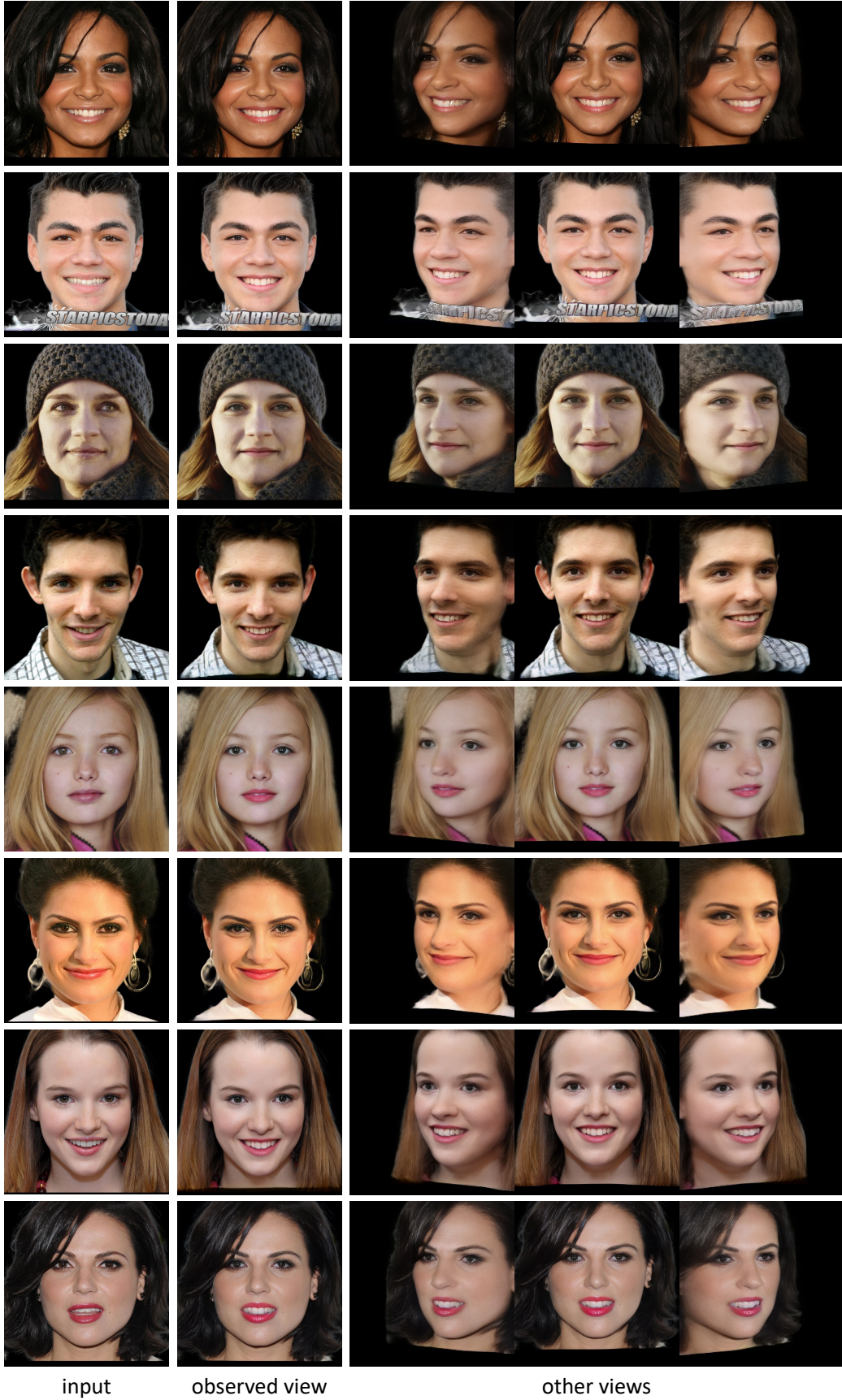
Figure 11: **Cross-identity reenactment on CelebA.**

input        observed view                other views

Figure 12: **3D reconstruction on CelebA.**

input       observed view       other views
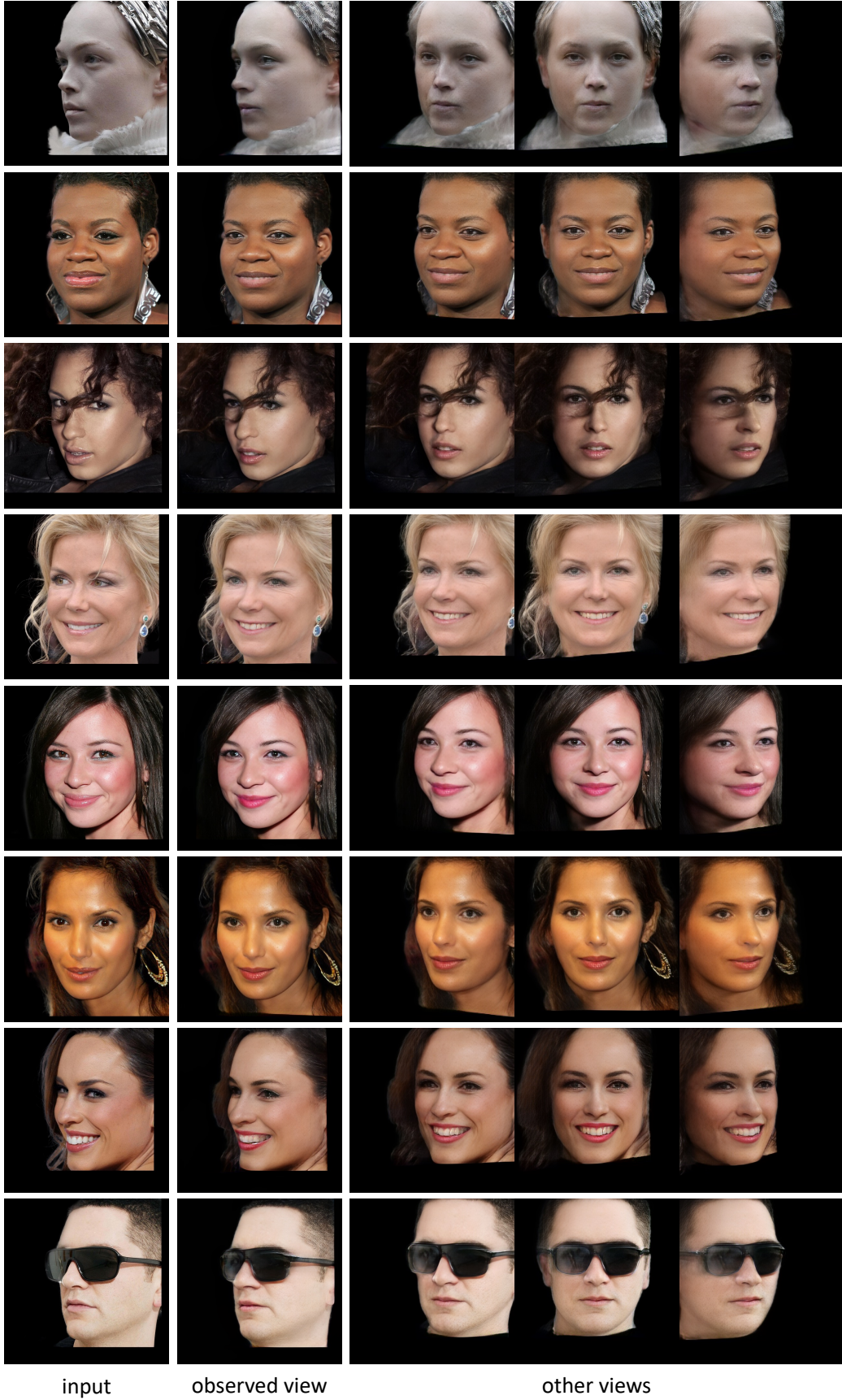
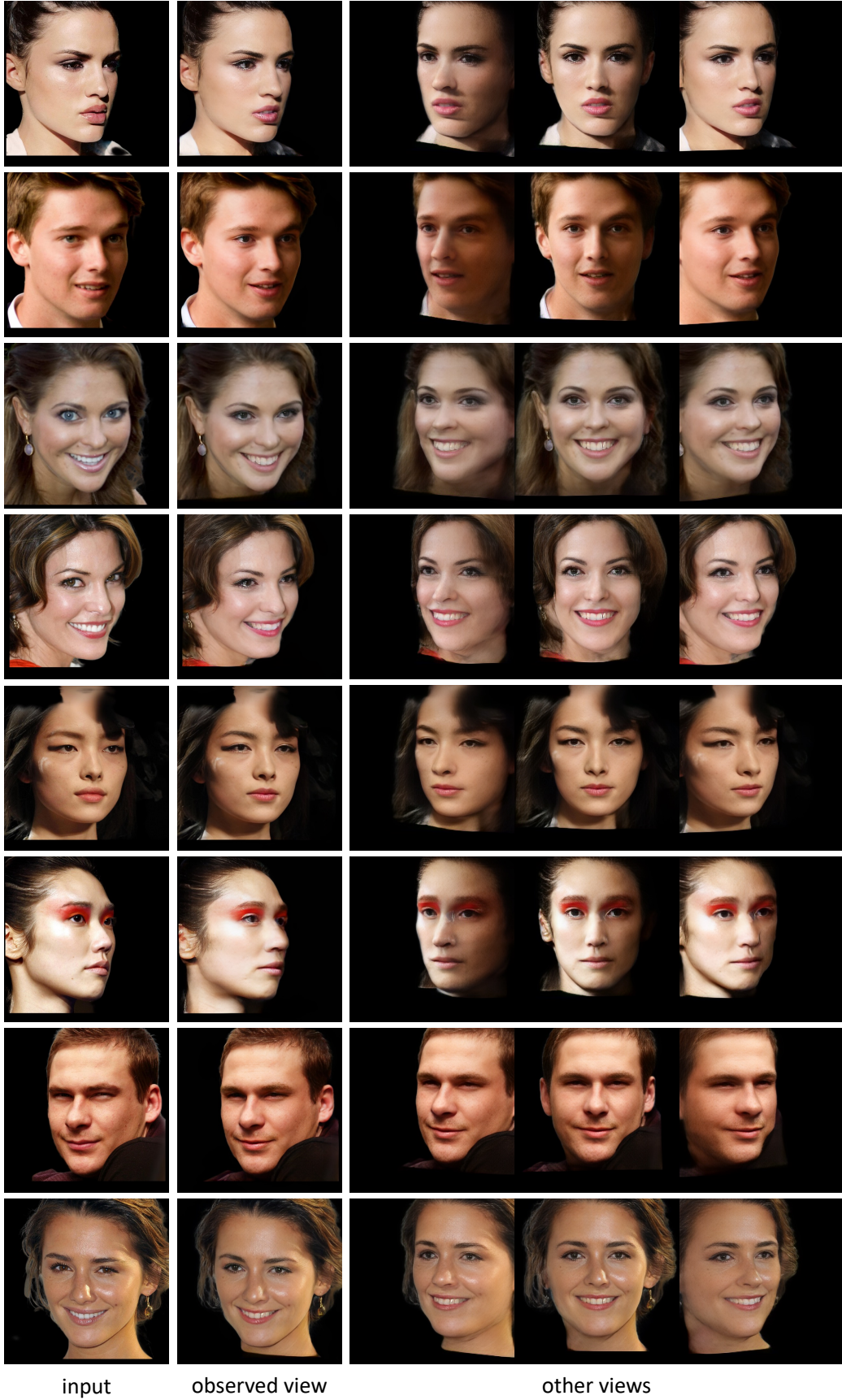Figure 13: **3D reconstruction on CelebA.**

input      observed view          other views

Figure 14: **3D reconstruction on CelebA.**

input        observed view        other views

Figure 15: **3D reconstruction on CelebA.**

Figure 16: **Cross-identity reenactment from HDTF to CelebA.**

(a) source  (b) target  (c) ours  (d) ROME  (e) StyleHeat  (f) Next3D-PTI



Figure 17: **Cross-identity reenactment from HDTF to CelebA.**

(a) source    (b) target    (c) ours    (d) ROME    (e) StyleHeat    (f) Next3D-PTI



Figure 18: **Cross-identity reenactment from HDTF to CelebA.**

25

(a) source  (b) target  (c) ours  (d) ROME  (e) StyleHeat  (f) Next3D-PTI

Figure 19: **Cross-identity reenactment from HDTF to CelebA.**

(a) source  (b) target  (c) ours  (d) OTAvatar  (e) ROME  (f) StyleHeat  (g) Next3D-PTI

Figure 20: **Same-identity reenactment on HDTF.**

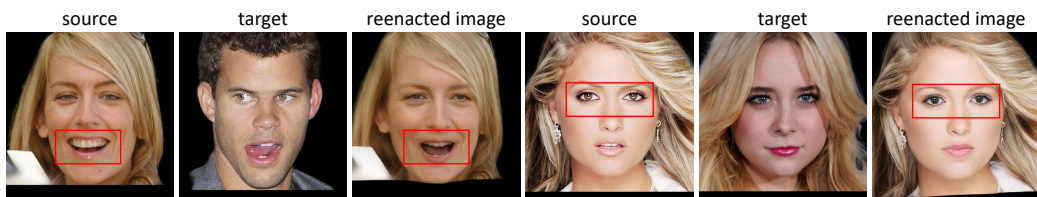Figure 21: **Cross-identity reenactment on HDTF.**



Figure 22: **Failure cases.**

# References

[1] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: a compact facial video forgery detection network. In *2018 IEEE international workshop on information forensics and security (WIFS)*, 2018. 10

[2] ShahRukh Athar, Zexiang Xu, Kalyan Sunkavalli, Eli Shechtman, and Zhixin Shu. Rignerf: Fully controllable neural 3d portraits. In *CVPR*, 2022. 2, 3

[3] Yunpeng Bai, Yanbo Fan, Xuan Wang, Yong Zhang, Jingxiang Sun, Chun Yuan, and Ying Shan. High-fidelity facial avatar reconstruction from monocular video with generative priors. *arXiv preprint arXiv:2211.15064*, 2022. 3

[4] Ziqian Bai, Feitong Tan, Zeng Huang, Kripasindhu Sarkar, Danhang Tang, Di Qiu, Abhimitra Meka, Ruofei Du, Mingsong Dou, Sergio Orts-Escolano, et al. Learning personalized high quality volumetric head avatars from monocular rgb videos. *arXiv preprint arXiv:2304.01436*, 2023. 2, 3

[5] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *ICCV*, 2021. 2, 3

[6] Ananta R Bhattarai, Matthias Nießner, and Artem Sevastopolsky. Triplanenet: An encoder for eg3d inversion. *arXiv preprint arXiv:2303.13497*, 2023. 3

[7] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. 2

[8] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *ICCV*, 2017. 8, 13

[9] 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, 4, 6, 8, 13

[10] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensoRF: Tensorial Radiance Fields. In *ECCV*, 2022. 2, 4

[11] Davide Cozzolino, Andreas Rössler, Justus Thies, Matthias Nießner, and Luisa Verdoliva. Id-reveal: Identity-aware deepfake video detection. In *CVPR*, 2021. 10

[12] Davide Cozzolino, Justus Thies, Andreas Rössler, Christian Riess, Matthias Nießner, and Luisa Verdoliva. Forensictransfer: Weakly-supervised domain adaptation for forgery detection. *arXiv preprint arXiv:1812.02510*, 2018. 10

[13] Radek Danecek, Michael J. Black, and Timo Bolkart. EMOCA: Emotion driven monocular face capture and animation. In *CVPR*, 2022. 1, 2

[14] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 8

[15] Yu Deng, Baoyuan Wang, and Heung-Yeung Shum. Learning detailed radiance manifolds for high-fidelity and 3d-consistent portrait synthesis from monocular image. *arXiv preprint arXiv:2211.13901*, 2022. 3, 4

[16] 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 *CVPRW*, 2019. 1, 2, 3, 4, 5, 8, 11

[17] Nikita Drobyshev, Jenya Chelishev, Taras Khakhulin, Aleksei Ivakhnenko, Victor Lempitsky, and Egor Zakharov. Megaportraits: One-shot megapixel neural head avatars. *Proceedings of the 30th ACM International Conference on Multimedia*, 2022. 1, 2

[18] Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)*, 2020. 2

[19] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3D face model from in-the-wild images. *SIGGRAPH*, 2021. 1, 2

[20] Anna Frühstück, Nikolaos Sarafianos, Yuanlu Xu, Peter Wonka, and Tony Tung. VIVE3D: Viewpoint-independent video editing using 3D-aware GANs. In *CVPR*, 2023. 3

[21] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, 2021. 2, 3

[22] Chen Gao, Yichang Shih, Wei-Sheng Lai, Chia-Kai Liang, and Jia-Bin Huang. Portrait neural radiance fields from a single image. *arXiv preprint arXiv:2012.05903*, 2020. 2, 3

[23] Xuan Gao, Chenglai Zhong, Jun Xiang, Yang Hong, Yudong Guo, and Juyong Zhang. Reconstructing personalized semantic facial nerf models from monocular video. *SIGGRAPH Asia*, 2022. 2, 3

[24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 1, 2, 6

[25] John C Gower. Generalized procrustes analysis. *Psychometrika*, 1975. 13

[26] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, 2021. 2, 3

[27] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *CVPR*, 2022. 3, 7, 8, 9, 13

[28] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 4

[29] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 6

[30] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 1, 6, 11

[31] Zhanghan Ke, Jiayu Sun, Kaican Li, Qiong Yan, and Rynson W.H. Lau. Modnet: Real-time trimap-free portrait matting via objective decomposition. In *AAAI*, 2022. 8

[32] Taras Khakhulin, Vanessa Sklyarova, Victor Lempitsky, and Egor Zakharov. Realistic one-shot mesh-based head avatars. In *ECCV*, 2022. 1, 3, 7, 8, 9, 12, 13

[33] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 8

[34] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nersemble: Multi-view radiance field reconstruction of human heads, 2023. 3

[35] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 7, 8, 9, 11, 12, 13

[36] Yeonkyeong Lee, Taeho Choi, Hyunsung Go, Hyunjoon Lee, Sunghyun Cho, and Junho Kim. Exp-gan: 3d-aware facial image generation with expression control. In *ACCV*, 2022. 3, 6

[37] Biwen Lei, Jianqiang Ren, Mengyang Feng, Miaomiao Cui, and Xuansong Xie. A hierarchical representation network for accurate and detailed face reconstruction from in-the-wild images, 2023. 2

[38] Ruilong Li, Karl Bladin, Yajie Zhao, Chinmay Chinara, Owen Ingraham, Pengda Xiang, Xinglei Ren, Pratusha Prasad, Bipin Kishore, Jun Xing, et al. Learning formation of physically-based face attributes. In *CVPR*, 2020. 2

[39] 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, (Proc. SIGGRAPH Asia)*, 2017. 2, 3

[40] Weichuang Li, Longhao Zhang, Dong Wang, Bin Zhao, Zhigang Wang, Mulin Chen, Bang Zhang, Zhongjian Wang, Liefeng Bo, and Xuelong Li. One-shot high-fidelity talking-head synthesis with deformable neural radiance field. In *CVPR*, 2023. 3

[41] Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. *PloS one*, 2018. 6

[42] Zhiyuan Ma, Xiangyu Zhu, Guojun Qi, Zhen Lei, and Lei Zhang. Otavatar: One-shot talking face avatar with controllable tri-plane rendering. *arXiv preprint arXiv:2303.14662*, 2023. 1, 3, 7, 8, 9, 12, 13

[43] 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*, 2020. 1, 2, 3

[44] Maryam Sadat Mirzaei, Kourosh Meshgi, Etienne Frigo, and Toyoaki Nishida. Animgan: A spatiotemporally-conditioned generative adversarial network for character animation. In *ICIP*, 2020. 3

[45] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 2022. 2, 3

[46] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 2, 3

[47] 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. *ACM Trans. Graph.*, 2021. 2, 3

[48] Frederic Ira Parke. *A parametric model for human faces.* The University of Utah, 1974. 2

[49] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 8

[50] P. Paysan, R. Knothe, B. Amberg, S. Romdhani, and T. Vetter. A 3d face model for pose and illumination invariant face recognition. *Proceedings of the 6th IEEE International Conference on Advanced Video and Signal based Surveillance (AVSS) for Security, Safety and Monitoring in Smart Environments*, 2009. 2, 14

[51] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *ICCV*, 2021. 1, 2

[52] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Trans. Graph.*, 2021. 3, 8, 13

[53] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. FaceForensics++: Learning to detect manipulated facial images. In *ICCV*, 2019. 10

[54] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *NeurIPS*, 2019. 1, 2

[55] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4

[56] 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. *arXiv preprint arXiv:2211.11208*, 2022. 3, 6, 8, 9, 12, 13

[57] Keqiang Sun, Shangzhe Wu, Zhaoyang Huang, Ning Zhang, Quan Wang, and HongSheng Li. Controllable 3d face synthesis with conditional generative occupancy fields. *arXiv preprint arXiv:2206.08361*, 2022. 3, 6

[58] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. *arXiv preprint arXiv:2303.14184*, 2023. 5

[59] Junshu Tang, Bo Zhang, Binxin Yang, Ting Zhang, Dong Chen, Lizhuang Ma, and Fang Wen. Explicitly controllable 3d-aware portrait generation. *arXiv preprint arXiv:2209.05434*, 2022. 3

[60] Kartik Teotia, Mallikarjun B R, Xingang Pan, Hyeongwoo Kim, Pablo Garrido, Mohamed Elgharib, and Christian Theobalt. Hq3davatar: High quality controllable 3d head avatar. In *arXiv*, 2023. 2, 3

[61] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 2, 3

[62] Alex Trevithick, Matthew Chan, Michael Stengel, Eric R Chan, Chao Liu, Zhiding Yu, Sameh Khamis, Manmohan Chandraker, Ravi Ramamoorthi, and Koki Nagano. Real-time radiance fields for single-image portrait view synthesis. *arXiv preprint arXiv:2305.02310*, 2023. 3

[63] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *CVPR*, 2021. 1, 2

[64] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *CVPR*, 2021. 1, 6, 13

[65] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *NeurIPS*, 2021. 4, 13

[66] Jiaxin Xie, Hao Ouyang, Jingtan Piao, Chenyang Lei, and Qifeng Chen. High-fidelity 3d gan inversion by pseudo-multi-view optimization. *arXiv preprint arXiv:2211.15662*, 2022. 3

[67] Eric Zhongcong Xu, Jianfeng Zhang, Junhao Liew, Wenqing Zhang, Song Bai, Jiashi Feng, and Mike Zheng Shou. Pv3d: A 3d generative model for portrait video generation. In *ICLR*, 2023. 3

[68] 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. *arXiv preprint arXiv:2303.15539*, 2023. 3

[69] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. *arXiv preprint arXiv:2201.08845*, 2022. 5

[70] Yiran Xu, Zhixin Shu, Cameron Smith, Jia-Bin Huang, and Seoung Wug Oh. In-n-out: Face video inversion and editing with volumetric decomposition. *arXiv preprint arXiv:2302.04871*, 2023. 3

[71] Fei Yang, Qian Zhang, Chi Zheng, and Guoping Qiu. In-the-wild facial expression recognition in extreme poses. In *International Conference on Graphic and Image Processing (ICGIP 2017)*, 2018. 13

[72] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. In *ECCV*, 2022. 1, 2, 7, 8, 9, 12, 13

[73] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *ECCV*, 2018. 13

[74] Ziyang Yuan, Yiming Zhu, Yu Li, Hongyu Liu, and Chun Yuan. Make encoder great again in 3d gan inversion through geometry and occlusion-aware encoding. *arXiv preprint arXiv:2303.12326*, 2023. 3, 4

[75] Egor Zakharov, Aleksei Ivakhnenko, Aliaksandra Shysheya, and Victor Lempitsky. Fast bi-layer neural synthesis of one-shot realistic head avatars. In *ECCV*, 2020. 1, 2

[76] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. Fdnerf: Few-shot dynamic neural radiance fields for face reconstruction and expression editing. In *SIGGRAPH Asia*, 2022. 3

[77] 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

[78] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. *arXiv preprint arXiv:2211.12194*, 2022. 1, 2

[79] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *CVPR*, 2021. 7, 8, 9, 12

[80] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *CVPR*, 2021. 10

[81] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I M Avatar: Implicit morphable head avatars from videos. In *CVPR*, 2022. 2, 3

[82] Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J. Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. In *CVPR*, 2023. 3

[83] Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and Chen Change Loy. CelebV-HQ: A large-scale video facial attributes dataset. In *ECCV*, 2022. 6

[84] Peiye Zhuang, Liqian Ma, Sanmi Koyejo, and Alexander Schwing. Controllable radiance fields for dynamic face synthesis. In *3DV*, 2022. 3

[85] Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. Mofanerf: Morphable facial neural radiance field. In *ECCV*, 2022. 3, 8