# MOST-GAN: 3D Morphable StyleGAN for Disentangled Face Image Manipulation

**Safa C. Medin**[1,2]    **Bernhard Egger**[2,3]    **Anoop Cherian**[1]    **Ye Wang**[1]

**Joshua B. Tenenbaum**[2]    **Xiaoming Liu**[4]    **Tim K. Marks**[1]

[1]Mitsubishi Electric Research Laboratories (MERL)
[2]Massachusetts Institute of Technology
[3]Friedrich-Alexander-University Erlangen-Nuremberg
[4]Michigan State University

## Abstract

Recent advances in generative adversarial networks (GANs) have led to remarkable achievements in face image synthesis. While methods that use style-based GANs can generate strikingly photorealistic face images, it is often difficult to control the characteristics of the generated faces in a meaningful and disentangled way. Prior approaches aim to achieve such semantic control and disentanglement within the latent space of a previously trained GAN. In contrast, we propose a framework that a priori models physical attributes of the face such as 3D shape, albedo, pose, and lighting explicitly, thus providing disentanglement by design. Our method, MOST-GAN, integrates the expressive power and photorealism of style-based GANs with the physical disentanglement and flexibility of nonlinear 3D morphable models, which we couple with a state-of-the-art 2D hair manipulation network. MOST-GAN achieves photorealistic manipulation of portrait images with fully disentangled 3D control over their physical attributes, enabling extreme manipulation of lighting, facial expression, and pose variations up to full profile view.

## Introduction

Changing certain attributes of a given portrait image, also referred to as *face image manipulation*, is a popular research topic that demonstrates the synergy between computer vision and computer graphics. Face image manipulation has a wide range of applications such as varying the illumination conditions to make a portrait image more appealing (Sun et al. 2019), changing the identity of a person to anonymize an image (Gafni, Wolf, and Taigman 2019), and exchanging the hairstyle in a virtual try-out setting (Tan et al. 2020). Two key factors make face image manipulation particularly challenging. First, the human visual system is sensitive to the smallest artifacts in synthesized face images, and careful handling of detail is therefore crucial to achieve photorealism. Second, faces are 3D objects with rich variations in shape, expression, and appearance, and inferring such 3D variations from 2D images is inherently an ill-posed problem.

StyleGAN2 (Karras et al. 2020) is currently one of the most advanced models for 2D image generation, reaching unprecedented quality and photorealism in synthesizing face images. At the same time, 3D face models, such as those based on 3D Morphable Models (3DMMs) (Blanz and Vetter 1999; Egger et al. 2020), are commonly used in recovering 3D

faces from 2D images; however, these reconstruction methods often lack photorealism (Tewari et al. 2017; Deng et al. 2019b). There are a few recent approaches that aim to combine the physically grounded modeling of 3DMMs with the synthesizing capabilities of style-based GANs (Tewari et al. 2020a,b). However, these approaches build on a fixed generative model, StyleGAN, and apply the explicit 3D model as a guiding tool to disentangle the learned StyleGAN latent space. As a result, these models cannot escape the data manifold characterized by a trained StyleGAN. Thus, while they provide some amount of control, they lack the generalization capabilities or physical disentanglement of 3D models, which limits their ability to synthesize large variations in the physical attributes of a face image.

In this work, we propose a nonlinear 3D face model that explicitly separates shape, albedo, lighting, and pose, which we refer to as *physical attributes*. Since we represent each of these attributes explicitly, we are able to control each of them independently, either within their learned latent spaces or by direct manipulation of their 3D physical realization. By processing each physical attribute separately, our novel real-image manipulation method achieves full disentanglement of these attributes. This is in sharp contrast to state-of-the-art (SOTA) methods such as Deng et al. (2020); Tewari et al. (2020a); Groueix et al. (2018), in which entanglement among different attributes is inevitable as they are all represented in one common latent space. Our model combines the photorealism of style-based GAN architectures with the generalization capabilities of 3DMMs, which allows for extrapolating beyond the variations present in the datasets. As a result, our method is able to manipulate faces to new poses, expressions, and illumination conditions that are not well represented in the training set. We also couple our 3D face model with a SOTA 2D hair model (Tan et al. 2020) to achieve a complete portrait image manipulation pipeline, allowing for joint face and hair processing. The contributions of this work include:

- We present a novel face image manipulation method, 3D **MO**rphable **ST**yle**GAN** (MOST-GAN), which by design achieves full disentanglement of shape, albedo, lighting, pose, and hair.
- We successfully combine the generalization capabilities of 3DMMs with the photorealism of style-based GANs, which enables us to synthesize novel 3D-grounded portrait
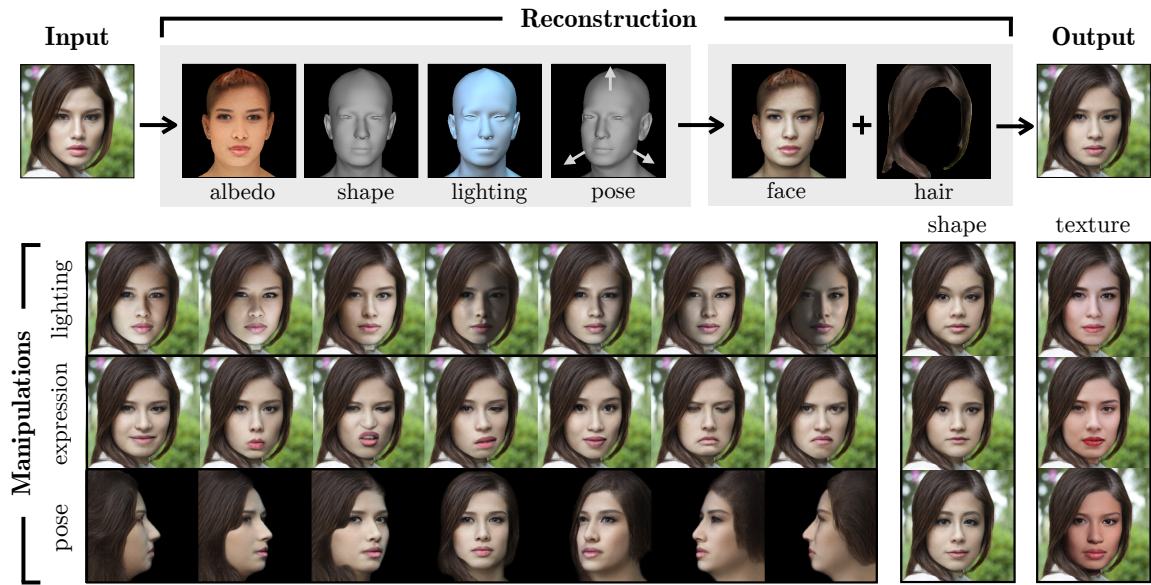
Figure 1: Fully disentangled, 3D controllable portrait image manipulation with MOST-GAN.

images with extreme variations that are rare or nonexistent in the training data.

- We develop a 3D-guided 2D hair manipulation algorithm, allowing for photorealistic and consistent hair styles and appearances over pose variations up to full profile views.

## Related Work

**Generative adversarial networks.** Generative adversarial networks (GANs) (Goodfellow et al. 2014) have set new standards in photorealistic image generation, with recent style-based methods StyleGAN (Karras, Laine, and Aila 2019) and StyleGAN2 (Karras et al. 2020) generating faces that are barely distinguishable from real photos. As conventional GANs learn only 2D representations, several works propose 3D GANs to achieve better understanding of the 3D world, via voxel-based (Choy et al. 2016; Wu et al. 2016, 2017; Zhu et al. 2018; Nguyen-Phuoc et al. 2019; Xie et al. 2019; Nguyen-Phuoc et al. 2020; Lunz et al. 2020) or mesh-based representations (Wang et al. 2018; Groueix et al. 2018; Pan et al. 2019). Recently, neural implicit representations have facilitated continuous 3D scene synthesis, including 3D faces (Schwarz et al. 2020; Chan et al. 2020). These methods, however, allow only limited control of facial pose. In another line of work, the 3D scene information is extracted from 2D GANs such as StyleGAN2 to manipulate 2D images in 3D (Shen and Zhou 2020; Härkönen et al. 2020) and recover explicit 3D shapes from images (Pan et al. 2020; Zhang et al. 2020). However, these methods do not employ strong shape priors such as 3DMMs, limiting their 3D manipulation capabilities. In contrast, we start from a 3D architecture while incorporating StyleGAN2 inside our network, which we train without using real 3D data.

**3D Morphable Models.** There is a classic line of research based on 3D Morphable Models (3DMMs) (Blanz and Vetter 1999; Egger et al. 2020) that aims for an object-specific 3D

model for faces based on high-quality 3D scans. Conventional linear 3DMMs such as the Basel Face Model (Paysan et al. 2009; Gerig et al. 2018) and FLAME (Li et al. 2017) typically suffer from a lack of expressiveness, due to their simplistic PCA-based texture and shape models and limited training data. To improve the representational power of 3DMMs, Tran and Liu (2018, 2019); Tran, Liu, and Liu (2019) proposed a nonlinear 3DMM that achieves better reconstruction quality than linear 3DMMs. Nonlinear models based on deep neural networks have also been used for realistic texture synthesis for various tasks (Saito et al. 2017; Slossberg, Shamai, and Kimmel 2018; Nagano et al. 2018). In this work, we build our face model as a nonlinear 3DMM based on the FLAME topology. Although the linear bases of FLAME do not yield photorealistic images, we use them to generate synthetic images for pretraining and to regularize our albedo reconstructions. In contrast to Tran and Liu (2018), we use separate encoders for different face attributes to foster further disentanglement among them, and employ StyleGAN2 for albedo synthesis, which generates images with better photorealism.

**3D Face Reconstruction.** A key application of 3DMMs is to reconstruct 3D faces from 2D images, with the objective to recover either face shape (Sanyal et al. 2019; Feng et al. 2020) or both shape and albedo. Methods that recover both shape and albedo have benefited from advancements in GANs, which enable higher quality and more realistic texture synthesis (Slossberg, Shamai, and Kimmel 2018; Gecer et al. 2019; Lattas et al. 2020). Among these approaches, GAN-FIT (Gecer et al. 2019) and AvatarMe (Lattas et al. 2020) obtain face reconstructions with high-frequency details, but they require large 3D datasets for training. Unlike those methods, ours does not rely on high-quality 3D data for photorealism—instead, we learn to generate detailed 3D face representations from 2D face images. Several other methods also recover 3D faces from only 2D images (Tewari et al. 2017; Deng et al. 2019b), although their reconstructions cannot be used
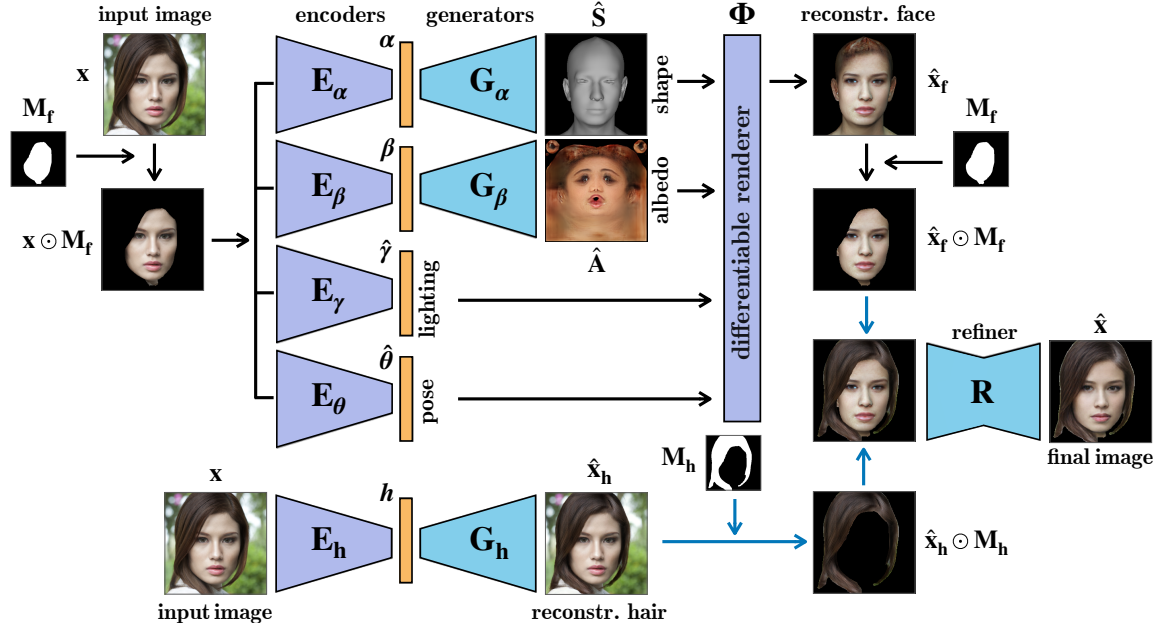
Figure 2: **Overview of our architecture.** Our model starts with a set of encoders for shape, albedo, lighting, pose, and hairstyle given an input image. To reconstruct the shape and albedo in their physical spaces, we use a convolutional generator for shape and a StyleGAN2 architecture for albedo. A hair generator reconstructs the hair in 2D. Reconstructed face and hair are finally fused and improved using a refiner network. All components are trained end-to-end, except for hair where we deploy a pretrained MichiGAN model (Tan et al. 2020).

for manipulating faces due to the lack of photorealism and missing details such as hair or teeth.

**Face image manipulation.** Recent research has aimed to combine 3DMMs with SOTA GANs to edit portrait images in a disentangled manner (Usman et al. 2019; Kowalski et al. 2020; Deng et al. 2020; Tewari et al. 2020a,b; Ghosh et al. 2020; Bühler et al. 2021; Piao et al. 2021). Among them, DiscoFaceGAN (Deng et al. 2020) promotes disentanglement between face attributes via contrastive learning, while StyleRig (Tewari et al. 2020a) couples a 3DMM with a pretrained 2D StyleGAN and manipulates images in the latent space of the StyleGAN. Since both methods rely on 2D generative networks, they are not able to handle extreme variations in 3D such as extreme lighting, facial expression, or pose. Furthermore, to manipulate real images, they must be embedded into the learned latent spaces via Image2StyleGAN (Abdal, Qin, and Wonka 2019), which hinders the quality of results. To circumvent such issues, Portrait Image Embedding (PIE) (Tewari et al. 2020b) introduces a novel optimization algorithm to embed real images into the latent space while preserving their photorealism. However, since both StyleRig and PIE are built on a pretrained 2D StyleGAN, the learned latent space limits them to variations that are well represented in the FFHQ dataset (Karras, Laine, and Aila 2019). Further, since these methods aim to disentangle their latent spaces *post hoc*, full disentanglement between physical attributes cannot be attained. In work concurrent to our research, GAR (Piao et al. 2021) proposes a realistic face reconstruction method that is used to manipulate portrait images, and VariTex (Bühler et al. 2021) introduces a variational texture generator to synthesize realistic face images while achieving control over

them. Both of these methods, however, provide only head pose and expression manipulation results, and similar to the other methods presented in this paragraph, they do not show results with large pose variations (larger than $45°$).

**Societal impacts.** We envision that our method could be used in numerous applications including creative uses in the entertainment sector, generation of realistic training data, or anonymization of public images. We have seen previous related methods misused to produce malicious content, such as fake news, and our method could enable face editing with larger variations. However, the research community is simultaneously creating methods to detect and mitigate such applications (Ciftci, Demir, and Yin 2020), and legal regulations to prohibit such misuse are under consideration.

## Methods

Our approach combines a statistical model of 3D faces with a style-based GAN, achieving a realistic and fully disentangled 3D model of faces. We achieve such disentanglement by individually processing each of the face's physical attributes and hair in the architecture, through separate encoders and decoders, as shown in Fig. 2. Such explicit control enables us to extrapolate beyond what is well represented in the training set, allowing for face synthesis in extreme poses, facial expressions, and lighting conditions.

### Problem Formulation

Our face image manipulation method relies on reconstructing accurate and photorealistic 3D faces from 2D images using the architecture shown in Fig. 2. Here, we assume that a portrait image can be decomposed into five different

attributes: four *physical attributes* (3D shape, albedo, lighting, and pose), and hair. Our face model employs a set of encoders $\{\mathbf{E}_\alpha, \mathbf{E}_\beta, \mathbf{E}_\gamma, \mathbf{E}_\theta\}$. Given a masked face image $\mathbf{x}' := \mathbf{x} \odot \mathbf{M_f}$, where $\mathbf{x}$ denotes the input image and $\mathbf{M_f}$ denotes its estimated face mask (Chen et al. 2017), the encoders $\mathbf{E}_\alpha$ and $\mathbf{E}_\beta$ extract a latent shape code $\alpha$ and albedo code $\beta$, while $\mathbf{E}_\gamma$ and $\mathbf{E}_\theta$ directly estimate the lighting parameters $\hat{\gamma}$ and pose parameters $\hat{\theta}$. To generate a face image, the shape and albedo codes are fed to a shape generator $\mathbf{G}_\alpha$ and albedo generator $\mathbf{G}_\beta$, respectively, to produce a 3D shape $\hat{\mathbf{S}}$ and albedo map $\hat{\mathbf{A}}$. Next, a differentiable renderer $\Phi$ renders the generated 3D model $\{\hat{\mathbf{S}}, \hat{\mathbf{A}}\}$ using the lighting and pose parameters $\{\hat{\gamma}, \hat{\theta}\}$ to produce the reconstructed face $\hat{\mathbf{x}_f}$: $\hat{\mathbf{x}_f} = \Phi(\hat{\mathbf{S}}, \hat{\mathbf{A}}, \hat{\gamma}, \hat{\theta})$.

Our hair model consists of an encoder $\mathbf{E_h}$ and a generator $\mathbf{G_h}$ to produce a portrait image with reconstructed hair $\hat{\mathbf{x}_h}$. Finally, the outputs of the face model and the hair model are combined using a face mask $\mathbf{M_f}$ and a hair mask $\mathbf{M_h}$, then passed through a refiner network $\mathbf{R}$ that produces the final image $\hat{\mathbf{x}}$. Formally, given a set of $N$ portrait images along with their face masks and hair masks $\{(\mathbf{x}^i, \mathbf{M_f}^i, \mathbf{M_h}^i)\}_{i=1}^N$, our objective is to solve the following optimization problem:

$$\underset{\{\mathbf{E}_\alpha, \mathbf{E}_\beta, \mathbf{E}_\gamma, \mathbf{E}_\theta, \mathbf{G}_\alpha, \mathbf{G}_\beta, \mathbf{R}\}}{\arg \min} \sum_{i=1}^N \left\| \mathbf{x}^i \odot (\mathbf{M_f}^i + \mathbf{M_h}^i) - \hat{\mathbf{x}}^i \right\|_1 \tag{1}$$

where each final image $\hat{\mathbf{x}} = \mathbf{R}(\hat{\mathbf{x}_f} \odot \mathbf{M_f} + \hat{\mathbf{x}_h} \odot \mathbf{M_h})$, with $\hat{\mathbf{x}_f} = \Phi(\mathbf{G}_\alpha(\mathbf{E}_\alpha(\mathbf{x}')), \mathbf{G}_\beta(\mathbf{E}_\beta(\mathbf{x}')), \mathbf{E}_\gamma(\mathbf{x}'), \mathbf{E}_\theta(\mathbf{x}'))$ and $\hat{\mathbf{x}_h} = \mathbf{G_h}(\mathbf{E_h}(\mathbf{x}))$. In later sections, we will show that adopting this objective enables us to edit portrait images in a fully disentangled manner while preserving their photorealism.

## Face Model

Our face model, demarcated in Fig. 2 by black connecting arrows, consists of four physical attribute encoders, two generators, and a differentiable renderer (Ravi et al. 2020). In the shape pipeline, the shape code $\alpha$ is input to a convolutional generator, $\mathbf{G}_\alpha$. The generated 3D shape, $\hat{\mathbf{S}}$, is composed of 3 channels in the UV-space that represent the 3D coordinates of vertices (Tran and Liu 2018) by their displacement from the FLAME mean head model. In parallel, the albedo code $\beta$ goes through a StyleGAN2 (Karras et al. 2020) generator $\mathbf{G}_\beta$ that outputs an RGB albedo map $\hat{\mathbf{A}}$ in the UV-space. Since most of the variations in face images are due to the variations in the albedo, generating albedo with a style-based architecture is a crucial step to achieve realism in the final output. Furthermore, in order to allow for more expressive latent spaces of shape and albedo, we let our model learn them without being constrained to the subspace defined by the original 3DMM. Finally, we represent the estimated lighting $\hat{\gamma}$ using a spherical harmonics parameterization with 3 bands (Ramamoorthi and Hanrahan 2001; Zhang and Samaras 2006), and our 6-DOF pose vector $\hat{\theta}$ includes 3 parameters for 3D rotation using the axis-angle representation and 3 parameters for 3D translation.

We divide our training process into two stages: 1) we pretrain our face model on synthetically generated faces; then 2) we generalize our model to real faces by training on real 2D images. The loss functions for each stage are introduced in the equations below and the subsequent explanations:

### Synthetic data Pretraining

$$L_{\text{image}}^{\text{syn}} = \|\mathbf{x} - \hat{\mathbf{x}_f}\|_2^2 \tag{2}$$

$$L_{\text{albedo}}^{\text{syn}} = \|\mathbf{A} - \hat{\mathbf{A}}\|_2^2 \tag{3}$$

$$L_{\text{shape}}^{\text{syn}} = \|\mathbf{w_s}^T (\mathbf{S} - \hat{\mathbf{S}})\|_2^2 \tag{4}$$

$$L_{\text{pose}}^{\text{syn}} = \|\theta - \hat{\theta}\|_2^2 \tag{5}$$

$$L_{\text{lighting}}^{\text{syn}} = \|\gamma - \hat{\gamma}\|_2^2 \tag{6}$$

$$L_{\text{reg}}^{\text{syn}} = \lambda_\alpha \|\alpha\|_2^2 + \lambda_\beta \|\beta\|_2^2 \tag{7}$$

### Real data Training

$$L_{\text{image}}^{\text{real}} = \|\mathbf{x} \odot \mathbf{M_f} - \hat{\mathbf{x}} \odot \mathbf{M_f}\|_2^2 \tag{8}$$

$$L_{\text{identity}}^{\text{real}} = 1 - \cos(f_{\text{id}}(\mathbf{x}), \, f_{\text{id}}(\hat{\mathbf{x}}')) \tag{9}$$

$$L_{\text{landmark}}^{\text{real}} = \|\mathbf{w_l}^T [f_{\text{lmk}}^{(1)}(\mathbf{x}) - f_{\text{lmk}}^{(2)}(\hat{\mathbf{S}})]\|_2^2 \tag{10}$$

$$L_{\text{albedo}}^{\text{real}} = \|(\mathbf{B^T B})^{-1} \mathbf{B^T} (\hat{\mathbf{A}} - \bar{\mathbf{A}})\|_2^2 \tag{11}$$

$$L_{\text{lighting}}^{\text{real}} = (\hat{\gamma} - \bar{\gamma})^T \Sigma^{-1} (\hat{\gamma} - \bar{\gamma}) \tag{12}$$

$$L_{\text{reg}}^{\text{real}} = \lambda_\alpha \|\alpha\|_2^2 + \lambda_\beta \|\beta\|_2^2 \tag{13}$$

**Pretraining on Synthetic Data.** The first stage is a pretraining step to allow our network to capture important characteristics of faces using strong supervision coming from a linear 3DMM. In this stage, we use the FLAME model to sample $80,000$ faces under an illumination and pose prior (Deng et al. 2020). We translate each face in 3D so that the rendered faces have the same 2D alignment as the FFHQ faces. Although these synthetic faces lack realism, they have ground truth values for the disentangled physical attributes albedo $\mathbf{A}$, shape $\mathbf{S}$, pose $\theta$, and lighting $\gamma$, which we use to guide pretraining. Our loss function for pretraining consists of three parts: reconstruction losses for the reconstructed face image (2) and for the four physical attributes (3)–(6); regularization for shape and albedo codes (7); and a non-saturating logistic GAN loss (Goodfellow 2016) to improve photorealism. In the shape reconstruction loss (4), we introduce a weighting term $\mathbf{w_s}$ to upweight vertices in regions surrounding salient facial features (e.g., eyes, eyebrows, mouth).

**Training on Real Data.** After pretraining, we train our model using the FFHQ face dataset (Karras, Laine, and Aila 2019), where for simplicity we eliminate the images with glasses. We obtain the face mask $\mathbf{M_f}$ for each image automatically using a semantic segmentation network (Or-El et al. 2020; Chen et al. 2017), then feed the masked 2D face images to the network. We train our face model in an end-to-end fashion, where we combine the loss functions in (8)–(13) with a non-saturating logistic GAN loss. Since we do not know the ground truth physical attributes for the real face images, we cannot apply any of the physical attribute reconstruction losses (3)–(6). The only reconstruction loss we apply is a pixelwise reconstruction loss for the masked faces (8). Defining the full reconstructed image as $\hat{\mathbf{x}}' := \mathbf{x} \odot (1 - \mathbf{M_f}) + \hat{\mathbf{x}} \odot \mathbf{M_f}$, we impose an identity loss (9), where $f_{\text{id}}(\cdot)$ denotes the feature vector extracted by the Arcface face recognition net-
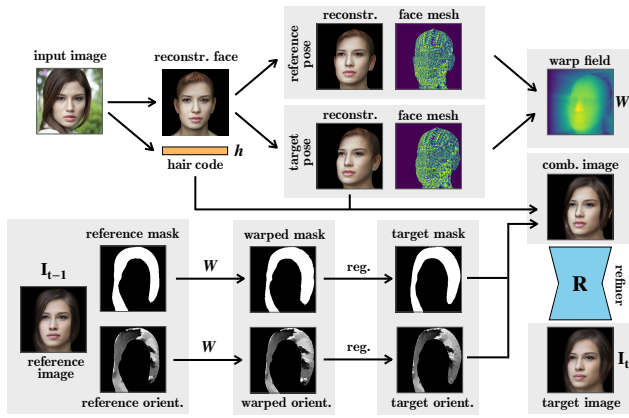
Figure 3: **One iteration of our hair manipulation algorithm**. Given a reference pose from the previous iteration and a target pose, we calculate a 2D warp field based on how 3D vertices move within the image plane. Given a reference image from the previous iteration $\mathbf{I_{t-1}}$ along with its reference mask and orientation, we use this warp field to warp the mask and the orientation, which we regularize to obtain the target mask and orientation. Next, we combine these with the hair appearance code obtained from the original input image and the reconstructed face reposed to the target pose, to obtain a novel portrait image $\mathbf{I_t}$. At the end, we feed this image through the refiner to obtain a photorealistic output. This algorithm is invoked sequentially starting from the original pose. The elements shown on gray backgrounds are updated in each iteration.

work (Deng et al. 2019a), and $\cos(\cdot, \cdot)$ denotes cosine similarity. Our landmark loss (10) measures the distance between the image-plane projections of the 3D facial landmark locations in the input image (estimated using (Bulat and Tzimiropoulos 2017)) and the corresponding locations in the reconstructed 3D shape model. The shape model vertices corresponding to specific facial landmarks are defined by the FLAME topology, and the weighting term $\mathbf{w_l}$ places more weight on important landmarks such as the lip outlines to keep our learned model faithful to the FLAME topology.

Since the decomposition of an input image into physical face properties is an ill-posed problem, there are ambiguities such as the relative contributions of color lighting intensities and surface albedo to the RGB appearance of a skin pixel. To help resolve this ambiguity, we introduce an albedo regularization loss (11) to minimize the projection of our reconstructed albedo into the FLAME model's albedo PCA space. Here, $\bar{\mathbf{A}}$ and $\mathbf{B}$ respectively represent the mean and basis vectors of the FLAME albedo model. To address the same ambiguity, we also include a lighting regularization loss (12), which minimizes the log-likelihood of the reconstructed lighting parameters $\hat{\gamma}$ under a multivariate Gaussian distribution over lighting conditions. To obtain that distribution, we sampled 50,000 lighting vectors using the prior provided by Deng et al. (2020) and calculate their sample mean $\bar{\gamma}$ and sample covariance $\Sigma$. As in pretraining, (13) regularizes the shape and albedo codes.

## Hair Model

Since hair has a more complex structure than faces, representing and manipulating hair in 3D is a very challenging

problem. This motivates us to manipulate hair in 2D, but to couple the hair generation process with our 3D face model. We build our hair model upon a state-of-the-art 2D model, MichiGAN (Tan et al. 2020), which disentangles hair shape, structure, and appearance by processing them separately and combines them with a backbone network. Here, shape refers to a 2D binary mask of the hair region, structure is represented as a 2D hair strand orientation map, and appearance refers to the global color and style of the hair which is encoded in a latent space. We incorporate a pretrained MichiGAN in our training pipeline, which we briefly represent as an encoder-decoder style model in Fig. 2. When we repose faces at inference time, we couple MichiGAN with our 3D face model to change the shape and structure of the hair without changing its appearance code.

**Coupling with Face Model.** Our 3D-guided hair manipulation algorithm is illustrated in Fig. 3. Since our face model reconstructs explicit 3D face shapes, we use these to reason about how the hair will move in 2D by calculating a 2D warp field (Li, Huang, and Loy 2019). We derive the 2D warp field based on the pose-induced movement of the 3D face vertices, then extrapolate the face's warp field to the rest of the image.

We use the warp field to warp the hair mask and the hair orientation map in 2D. Since this process can introduce warp artifacts, however, we regularize the warped masks by projecting them onto a PCA basis calculated from a dataset of binary hair masks of portrait images. In addition to obtaining hair masks from the FFHQ dataset (Karras, Laine, and Aila 2019), we extract hair masks from the USC HairSalon database (Hu et al. 2015) by rendering that dataset's 3D hair models with faces in extreme poses to allow for accurate and consistent hair masks under large pose variations. The orientation map, on the other hand, is regularized as part of the MichiGAN framework, which outputs a map that is consistent with the warped map and aligned with the regularized hair mask. Finally, the reconstructed face in the target pose, hair appearance code, hair mask, and hair orientation map are combined by the MichiGAN pipeline to produce the reposed portrait image, which is then processed with the refiner (described below). For large pose variations, we invoke this algorithm sequentially by going from reference pose to target pose in multiple steps, and we regularize the warped masks and orientation maps at each step. For more details, please see the supplementary material.

## Refinement

Although our combined model's rendered 3D face reconstructions and 2D hair reconstructions closely resemble the original images, there is still a small realism gap that needs to be filled. In particular, since we regularize the reconstructed albedos using the FLAME albedo space, the reconstructions do not exhibit sufficient variation in the eye regions, and they lack certain details such as eyelashes, facial hair, teeth, and accessories, which are not modeled by the FLAME mesh template. Furthermore, since face and hair are processed separately, some reconstructions have blending issues between the face and the hair. To address these issues, we utilize a refiner network, which closes the realism gap between the
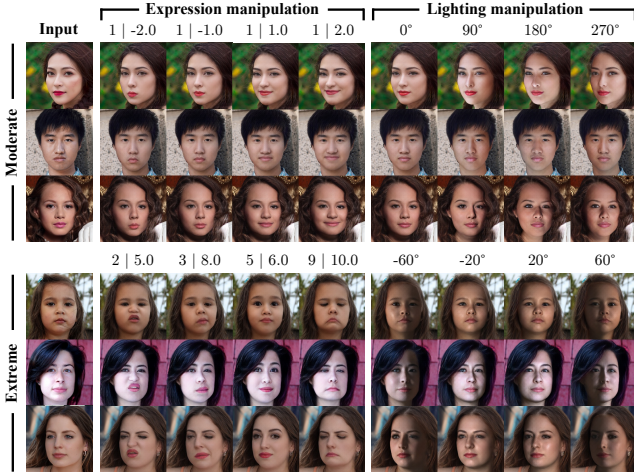
Figure 4: Expression and lighting manipulation results. **Expression manipulation (*left*).** We illustrate both moderate (*top*) and extreme (*bottom*) expression variations. The two numbers above each column indicate which FLAME expression eigenvector is used and by how many standard deviations it is scaled. **Lighting manipulation (*right*).** *Top:* For moderate variation, we rotate the reconstructed lighting around the camera axis by the angle above each column. *Bottom:* For extreme lighting variation, we render the reconstructed 3D model using a point light source and Phong shading model.

reconstructions and the original images while making only a minimal change to the reconstructions. We employ a U-Net (Ronneberger, Fischer, and Brox 2015) that takes in an image combining the reconstructed face and hair and outputs a more realistic portrait image, as shown in Fig. 2.

After freezing the weights of the rest of the model, we train the refiner with pairs of original images from the dataset and reconstructed images. For the refiner, we combine the same adversarial loss and identity loss (9) described above with a reconstruction loss based on the VGG-16 perceptual loss (Simonyan and Zisserman 2014; Johnson, Alahi, and Fei-Fei 2016), promoting better reconstruction quality for hair.

## Experiments and Results

In our experiments, we manipulate portrait images with respect to several physical attributes and compare them with a SOTA real-image manipulation method, PIE (Tewari et al. 2020b). Besides providing qualitative comparisons of the two methods, we also quantitatively compare the performance of our pose editing algorithm by employing a head pose estimator (Ruiz, Chong, and Rehg 2018) to measure the error between the desired and estimated head poses.

Because MOST-GAN generates a full 3D model, we can manipulate physical attributes beyond the distribution of the training set. We can also modify the face in ways not anticipated during training, such as relighting faces using a different lighting and shading model.

**Expression and lighting manipulation.** We illustrate our facial expression and lighting manipulation results in Fig. 4. To edit facial expression (left), we choose an eigenvector from the FLAME expression basis and multiply it by a constant factor to obtain an offset, which we add to the vertex
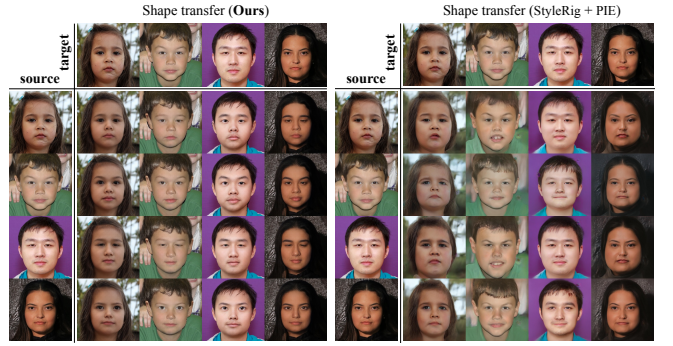


Figure 5: Shape transfer comparison. We transfer the 3D shape of each source image to each target image while keeping everything else unchanged. Our results (*left*) demonstrate more accurate shape transfer and much better disentanglement between shape and other attributes (e.g., albedo, pose, and hair) than the combination of StyleRig (Tewari et al. 2020a) and PIE (Tewari et al. 2020b) (*right*).

locations in our model's reconstructed 3D shape. In the moderate examples (top left), we use the first eigenvector to add smile/frown variations up to $\pm 2$ standard deviations. In the extreme examples (bottom left), we scale 5 different expression eigenvectors by up to 10 standard deviations. For lighting manipulation (right), the moderate edits (top right) rotate the reconstructed lighting around the camera axis (axis perpendicular to the image plane). For extreme lighting variations (bottom right), we employ a point light source and Phong shading model, where we rotate the light source horizontally around the vertical axis and can introduce any desired amount of specularity to the face albedo. The results demonstrate that our method easily handles extreme expressions and lighting conditions that are not well-represented in the training set and can use lighting and shading models not used in training. (We show more examples in the supplementary material.)

**Shape transfer.** Our model achieves superior disentanglement of physical attributes such as shape and albedo by design, by modeling them separately and explicitly. This disentanglement is illustrated by the shape transfer results in Fig. 5, where we transfer the 3D face shape of a source image to a target image. Our results show that our method (*left*) is able to transfer the face shapes accurately, while maintaining photorealism and keeping the albedo, lighting, and hair unchanged. This is in contrast to the shape transfer results by the previous state of the art (*right*, a combination of StyleRig (Tewari et al. 2020a) and PIE (Tewari et al. 2020b)), where for a given source shape, the transfer results have varying face shapes with noticeable differences in expressions. When the source and target images are identical (images on the diagonal), our method produces the original reconstruction by design, whereas PIE + StyleRig struggles to maintain the original identity. Our method can also transfer albedo alone, transfer multiple physical attributes (such as albedo and shape) simultaneously, and smoothly interpolate between different shapes and albedos in the latent space continuously. (See the supplementary material for examples.)

**Pose manipulation.** In Fig. 6, we compare our pose manipulation results (odd rows) to PIE (Tewari et al. 2020b) (even rows). To edit the pose of a given portrait image, we

Table 1: Mean absolute errors between the desired and estimated head poses in degrees, on 100 random images from our test set. Our method's average head pose error is significantly smaller than that of PIE (Tewari et al. 2020b), indicating our method's superior pose disentanglement.

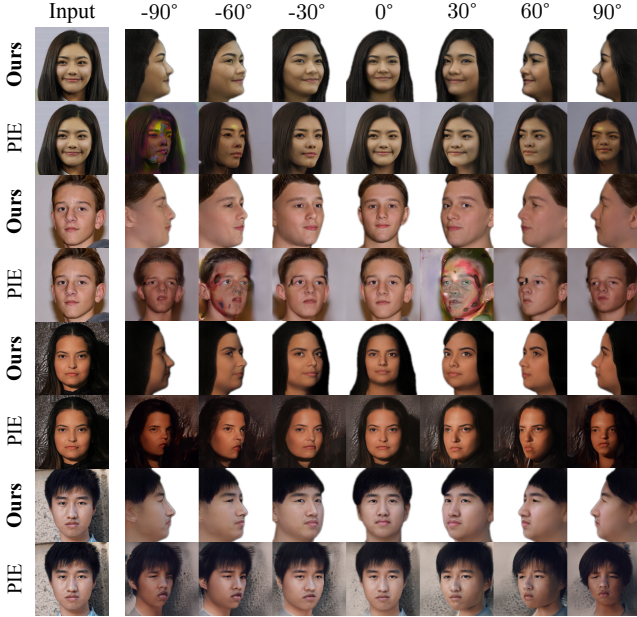| | -90° | -75° | -60° | -45° | -30° | -15° | 0° | 15° | 30° | 45° | 60° | 75° | 90° |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Ours** | **13.2** | **6.0** | **3.5** | **4.5** | **4.7** | **4.8** | **2.3** | **2.6** | **5.0** | **6.0** | **3.9** | **2.7** | **6.7** |
| PIE | 61.7 | 38.0 | 22.1 | 15.8 | 10.8 | 7.6 | 3.5 | 6.2 | 13.1 | 22.5 | 35.5 | 54.8 | 81.0 |



Figure 6: Pose manipulation results. From an input portrait image, our method accurately rotates the reconstructed 3D face all the way to profile pose. (Faces in more extreme poses appear larger due to the FFHQ alignment.) In contrast, PIE (Tewari et al. 2020b) struggles to maintain photorealism and cannot achieve large rotations.

rotate the reconstructed faces in 3D and warp the hair in 2D using our 3D-guided hair manipulation algorithm described in the Hair Model section. The results show that our method is able to rotate portrait images all the way to profile pose while keeping the identity, expression, and illumination conditions unchanged. For the 0° pose, PIE (Tewari et al. 2020b) is slightly better at reconstructing the original identity. However, PIE relies on a costly optimization over the latent space of a pretrained GAN, whereas our method reconstructs 3D faces at interactive framerates (30 fps) using our encoder-decoder style architecture. Furthermore, PIE cannot handle extreme rotations that were not well represented in the StyleGAN training set, yielding unrealistic artifacts and an inability to achieve larger desired (target) poses. To quantify the latter, we calculate the mean absolute error between the desired and achieved head poses using a head pose estimation network (Ruiz, Chong, and Rehg 2018). In particular, we randomly sampled 100 images from our test dataset, reposed them in a range of yaw angles using our method and PIE, and calculated the average absolute pose error of each method. The results, in Table 1, show that our method yields much more accurate pose manipulation at all pose angles.

**Limitations**. Since we disentangle hair from the physical attributes by design, changing the lighting conditions has a limited effect on the hair, and that effect is achieved by the refiner. Since the hair appearance is strongly dependent on the head pose and lighting conditions, this issue could be addressed by coupling the pose and lighting with the hair model at training time. Also, since the reconstruction quality of hair is heavily influenced by the hair orientation map in the MichiGAN framework, achieving consistency of orientation maps over large pose variation is crucial to render photorealistic hair for reposed images. Currently, however, we handle dis-occlusions of the hair by warping the orientation maps in 2D, which sometimes yields inconsistent orientations (and thus unrealistic hair rendering) after large pose changes. In addition, our model tends to attribute skin color mostly to the lighting component, which is due to the fact that samples from the FLAME albedo basis, which we use to regularize our albedo reconstructions, do not exhibit much variation in skin tone. Finally, our face model tends to yield smooth 3D shape reconstructions, sometimes attributing fine shape details such as wrinkles on the face to the albedo instead of the shape. We believe that this is related to our shape generator following a convolutional architecture, which promotes local consistency between neighboring vertices of the face mesh.

## Conclusion

In this work we introduce MOST-GAN, a novel framework for manipulating face images in a 3D controllable and fully disentangled way. We achieve this by combining the physically-grounded modeling of 3DMMs with the expressive power of style-based GANs. We employ an encoder-decoder style architecture built on a 3DMM template, where we represent 3D shape, albedo, pose, and lighting independently by design. By coupling our 3D face model with a state-of-the-art 2D hair model, we develop a full portrait image manipulation pipeline. Unlike state-of-the-art methods, which require costly optimizations before manipulating real images, our method enables efficient image manipulation at inference time. Our results demonstrate the ability of our method to photorealistically manipulate 3D shape, albedo, pose, and lighting of face images, facilitating larger variations compared to state-of-the-art methods, and achieving better disentanglement in face image manipulation tasks.

## Acknowledgements

# References

Abdal, R.; Qin, Y.; and Wonka, P. 2019. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4432–4441.

Blanz, V.; and Vetter, T. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 187–194.

Bühler, M. C.; Meka, A.; Li, G.; Beeler, T.; and Hilliges, O. 2021. VariTex: Variational Neural Face Textures. *arXiv preprint arXiv:2104.05988*.

Bulat, A.; and Tzimiropoulos, G. 2017. How far are we from solving the 2D & 3D Face Alignment problem? (and a dataset of 230,000 3D facial landmarks). In *International Conference on Computer Vision*.

Chan, E.; Monteiro, M.; Kellnhofer, P.; Wu, J.; and Wetzstein, G. 2020. pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis. In *arXiv*.

Chen, L.-C.; Papandreou, G.; Schroff, F.; and Adam, H. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*.

Choy, C. B.; Xu, D.; Gwak, J.; Chen, K.; and Savarese, S. 2016. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *European conference on computer vision*, 628–644. Springer.

Ciftci, U. A.; Demir, I.; and Yin, L. 2020. Fakecatcher: Detection of synthetic portrait videos using biological signals. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.

Deng, J.; Guo, J.; Niannan, X.; and Zafeiriou, S. 2019a. ArcFace: Additive Angular Margin Loss for Deep Face Recognition. In *CVPR*.

Deng, Y.; Yang, J.; Chen, D.; Wen, F.; and Tong, X. 2020. Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning. In *IEEE Computer Vision and Pattern Recognition*.

Deng, Y.; Yang, J.; Xu, S.; Chen, D.; Jia, Y.; and Tong, X. 2019b. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 0–0.

Egger, B.; Smith, W. A.; Tewari, A.; Wuhrer, S.; Zollhoefer, M.; Beeler, T.; Bernard, F.; Bolkart, T.; Kortylewski, A.; Romdhani, S.; et al. 2020. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)*, 39(5): 1–38.

Feng, Y.; Feng, H.; Black, M. J.; and Bolkart, T. 2020. Learning an Animatable Detailed 3D Face Model from In-The-Wild Images. *arXiv preprint arXiv:2012.04012*.

Gafni, O.; Wolf, L.; and Taigman, Y. 2019. Live face de-identification in video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9378–9387.

Gecer, B.; Ploumpis, S.; Kotsia, I.; and Zafeiriou, S. 2019. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1155–1164.

Gerig, T.; Morel-Forster, A.; Blumer, C.; Egger, B.; Luthi, M.; Schönborn, S.; and Vetter, T. 2018. Morphable face models-an open framework. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, 75–82. IEEE.

Ghosh, P.; Gupta, P. S.; Uziel, R.; Ranjan, A.; Black, M.; and Bolkart, T. 2020. Gif: Generative interpretable faces. *arXiv preprint arXiv:2009.00149*.

Goodfellow, I. 2016. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*.

Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*.

Groueix, T.; Fisher, M.; Kim, V. G.; Russell, B. C.; and Aubry, M. 2018. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 216–224.

Härkönen, E.; Hertzmann, A.; Lehtinen, J.; and Paris, S. 2020. Ganspace: Discovering interpretable gan controls. *arXiv preprint arXiv:2004.02546*.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.

Hore, A.; and Ziou, D. 2010. Image quality metrics: PSNR vs. SSIM. In *2010 20th international conference on pattern recognition*, 2366–2369. IEEE.

Hu, L.; Ma, C.; Luo, L.; and Li, H. 2015. Single-view hair modeling using a hairstyle database. *ACM Transactions on Graphics (ToG)*, 34(4): 1–9.

Johnson, J.; Alahi, A.; and Fei-Fei, L. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, 694–711. Springer.

Karras, T.; Laine, S.; and Aila, T. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4401–4410.

Karras, T.; Laine, S.; Aittala, M.; Hellsten, J.; Lehtinen, J.; and Aila, T. 2020. Analyzing and Improving the Image Quality of StyleGAN. In *Proc. CVPR*.

Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Kowalski, M.; Garbin, S. J.; Estellers, V.; Baltrušaitis, T.; Johnson, M.; and Shotton, J. 2020. Config: Controllable neural face image generation. *arXiv preprint arXiv:2005.02671*.

Lattas, A.; Moschoglou, S.; Gecer, B.; Ploumpis, S.; Triantafyllou, V.; Ghosh, A.; and Zafeiriou, S. 2020. AvatarMe: Realistically Renderable 3D Facial Reconstruction" In-the-Wild". In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 760–769.

Li, T.; Bolkart, T.; Black, M. J.; Li, H.; and Romero, J. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6): 194:1–194:17.

Li, Y.; Huang, C.; and Loy, C. C. 2019. Dense intrinsic appearance flow for human pose transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3693–3702.

Lunz, S.; Li, Y.; Fitzgibbon, A.; and Kushman, N. 2020. Inverse graphics gan: Learning to generate 3d shapes from unstructured 2d data. *arXiv preprint arXiv:2002.12674*.

Nagano, K.; Seo, J.; Xing, J.; Wei, L.; Li, Z.; Saito, S.; Agarwal, A.; Fursund, J.; and Li, H. 2018. paGAN: real-time avatars using dynamic textures. *ACM Transactions on Graphics (TOG)*, 37(6): 1–12.

Nguyen-Phuoc, T.; Li, C.; Theis, L.; Richardt, C.; and Yang, Y.-L. 2019. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7588–7597.

Nguyen-Phuoc, T.; Richardt, C.; Mai, L.; Yang, Y.-L.; and Mitra, N. 2020. BlockGAN: Learning 3D object-aware scene representations from unlabelled images. *arXiv preprint arXiv:2002.08988*.

Or-El, R.; Sengupta, S.; Fried, O.; Shechtman, E.; and Kemelmacher-Shlizerman, I. 2020. Lifespan Age Transformation Synthesis. In *European Conference on Computer Vision*, 739–755. Springer.

Pan, J.; Han, X.; Chen, W.; Tang, J.; and Jia, K. 2019. Deep mesh reconstruction from single rgb images via topology modification networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9964–9973.

Pan, X.; Dai, B.; Liu, Z.; Loy, C. C.; and Luo, P. 2020. Do 2d gans know 3d shape? unsupervised 3d shape reconstruction from 2d image gans. *arXiv preprint arXiv:2011.00844*.

Paysan, P.; Knothe, R.; Amberg, B.; Romdhani, S.; and Vetter, T. 2009. A 3D face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*, 296–301. Ieee.

Piao, J.; Sun, K.; Lin, K.; and Li, H. 2021. Inverting Generative Adversarial Renderer for Face Reconstruction. arXiv:2105.02431.

Ramamoorthi, R.; and Hanrahan, P. 2001. An efficient representation for irradiance environment maps. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 497–500.

Ravi, N.; Reizenstein, J.; Novotny, D.; Gordon, T.; Lo, W.-Y.; Johnson, J.; and Gkioxari, G. 2020. Accelerating 3D Deep Learning with PyTorch3D. *arXiv:2007.08501*.

Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 234–241. Springer.

Ruiz, N.; Chong, E.; and Rehg, J. M. 2018. Fine-grained head pose estimation without keypoints. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2074–2083.

Saito, S.; Wei, L.; Hu, L.; Nagano, K.; and Li, H. 2017. Photorealistic facial texture inference using deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5144–5153.

Sanyal, S.; Bolkart, T.; Feng, H.; and Black, M. J. 2019. Learning to regress 3D face shape and expression from an image without 3D supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7763–7772.

Schwarz, K.; Liao, Y.; Niemeyer, M.; and Geiger, A. 2020. Graf: Generative radiance fields for 3d-aware image synthesis. *arXiv preprint arXiv:2007.02442*.

Shen, Y.; and Zhou, B. 2020. Closed-form factorization of latent semantics in gans. *arXiv preprint arXiv:2007.06600*.

Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.

Slossberg, R.; Shamai, G.; and Kimmel, R. 2018. High quality facial surface and texture synthesis via generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 0–0.

Sun, T.; Barron, J. T.; Tsai, Y.-T.; Xu, Z.; Yu, X.; Fyffe, G.; Rhemann, C.; Busch, J.; Debevec, P. E.; and Ramamoorthi, R. 2019. Single image portrait relighting. *ACM Trans. Graph.*, 38(4): 79–1.

Tan, Z.; Chai, M.; Chen, D.; Liao, J.; Chu, Q.; Yuan, L.; Tulyakov, S.; and Yu, N. 2020. MichiGAN: Multi-Input-Conditioned Hair Image Generation for Portrait Editing. *ACM Transactions on Graphics (TOG)*, 39(4): 1–13.

Tewari, A.; Elgharib, M.; Bharaj, G.; Bernard, F.; Seidel, H.-P.; Pérez, P.; Zollhofer, M.; and Theobalt, C. 2020a. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6142–6151.

Tewari, A.; Elgharib, M.; BR, M.; Bernard, F.; Seidel, H.-P.; Pérez, P.; Zöllhofer, M.; and Theobalt, C. 2020b. PIE: Portrait Image Embedding for Semantic Control. *ACM Transactions on Graphics (Proceedings SIGGRAPH Asia)*, 39(6).

Tewari, A.; Zollhofer, M.; Kim, H.; Garrido, P.; Bernard, F.; Perez, P.; and Theobalt, C. 2017. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 1274–1283.

Tran, L.; Liu, F.; and Liu, X. 2019. Towards high-fidelity nonlinear 3D face morphable model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1126–1135.

Tran, L.; and Liu, X. 2018. Nonlinear 3d face morphable model. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7346–7355.

Tran, L.; and Liu, X. 2019. On learning 3d face morphable model from in-the-wild images. *IEEE transactions on pattern analysis and machine intelligence*, 43(1): 157–171.

Usman, B.; Dufour, N.; Saenko, K.; and Bregler, C. 2019. Puppetgan: Cross-domain image manipulation by demonstration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9450–9458.

Wang, N.; Zhang, Y.; Li, Z.; Fu, Y.; Liu, W.; and Jiang, Y.-G. 2018. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 52–67.

Wu, J.; Wang, Y.; Xue, T.; Sun, X.; Freeman, W. T.; and Tenenbaum, J. B. 2017. Marrnet: 3d shape reconstruction via 2.5 d sketches. *arXiv preprint arXiv:1711.03129*.

Wu, J.; Zhang, C.; Xue, T.; Freeman, W. T.; and Tenenbaum, J. B. 2016. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *arXiv preprint arXiv:1610.07584*.

Xie, H.; Yao, H.; Sun, X.; Zhou, S.; and Zhang, S. 2019. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2690–2698.

Zhang, L.; and Samaras, D. 2006. Face recognition from a single training image under arbitrary unknown lighting using spherical harmonics. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(3): 351–363.

Zhang, Y.; Chen, W.; Ling, H.; Gao, J.; Zhang, Y.; Torralba, A.; and Fidler, S. 2020. Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. *arXiv preprint arXiv:2010.09125*.

Zhu, J.-Y.; Zhang, Z.; Zhang, C.; Wu, J.; Torralba, A.; Tenenbaum, J. B.; and Freeman, W. T. 2018. Visual object networks: Image generation with disentangled 3d representation. *arXiv preprint arXiv:1812.02725*.

# Supplementary Material

## Architecture Details

**Encoders.** For encoders $\{\mathbf{E}_{\boldsymbol{\alpha}}, \mathbf{E}_{\boldsymbol{\beta}}, \mathbf{E}_{\boldsymbol{\gamma}}, \mathbf{E}_{\boldsymbol{\theta}}\}$, we employ the ResNet-18 architecture (He et al. 2016) (starting from the ImageNet (Deng et al. 2009) pretrained weights) where we change the final layer to reflect the dimensionality of each latent representation: $\boldsymbol{\alpha} \in \mathbb{R}^{150}$, $\boldsymbol{\beta} \in \mathbb{R}^{200}$, $\hat{\boldsymbol{\gamma}} \in \mathbb{R}^{27}$, and $\hat{\boldsymbol{\theta}} \in \mathbb{R}^6$.

**Generators.** For the albedo generator $\mathbf{G}_{\boldsymbol{\beta}}$, we employ the original StyleGAN2 (Karras et al. 2020) architecture up to the $256 \times 256$ layer. (We omit the final layers with higher resolutions.) For the shape generator $\mathbf{G}_{\boldsymbol{\alpha}}$, we use the architecture shown in Table 2. The output of this network is a UV-representation of shape from which we sample points corresponding to the UV-coordinates of each vertex in the FLAME (Li et al. 2017) topology. Then, we add these as an offset to the FLAME mean shape to obtain a 3D shape in Euclidean space.

**Refiner.** For the refiner, we employ a U-Net (Ronneberger, Fischer, and Brox 2015) with 5 convolutional layers followed by 5 transpose convolutional layers with skip connections. We provide the U-Net architecture details in Table 3.

## Training Details

Throughout this section, we let $L_G$ denote the *generator loss*, which is minimized to fool a discriminator that is trained adversarially to distinguish between generated face images and real face images (Goodfellow et al. 2014).

**Pretraining on Synthetic Data.** Our disentangled face model allows for pretraining each of the physical attributes separately. We carry out pretraining in three independent phases: albedo-only, lighting-only, and shape & pose jointly. Using the notation we have introduced before, we minimize the following loss functions for the three phases:

$$\text{shape \& pose:} \quad 0.1 L_{\text{image}}^{\text{syn}} + 1000 L_{\text{shape}}^{\text{syn}}$$
$$+ 100 L_{\text{pose}}^{\text{syn}} + 1.0 \|\boldsymbol{\alpha}\|_2^2$$
$$\text{albedo-only:} \quad L_G + 10 L_{\text{image}}^{\text{syn}} + 100 L_{\text{albedo}}^{\text{syn}} + 1.0 \|\boldsymbol{\beta}\|_2^2$$
$$\text{lighting-only:} \quad 10 L_{\text{image}}^{\text{syn}} + 100 L_{\text{light}}^{\text{syn}}$$

For each of these phases, we set the batch size to 16 and use the Adam optimizer (Kingma and Ba 2014). $\mathbf{E}_{\boldsymbol{\alpha}}$, $\mathbf{E}_{\boldsymbol{\beta}}$, $\mathbf{E}_{\boldsymbol{\gamma}}$, $\mathbf{E}_{\boldsymbol{\theta}}$, and $\mathbf{G}_{\boldsymbol{\alpha}}$ are all trained with a learning rate of 0.0001. During the albedo-only phase, we alternate optimization steps between training the albedo generator $\mathbf{G}_{\boldsymbol{\beta}}$ and the image discriminator (both with learning rate 0.002).

**Training on Real Data.** In the training stage, we split the FFHQ dataset (Karras, Laine, and Aila 2019) into train and test sets with $90\% - 10\%$ split, using the first 63,000 face images for training and the last 7,000 images for testing. During training, we optimize over all networks in our face model ($\mathbf{E}_{\boldsymbol{\alpha}}$, $\mathbf{E}_{\boldsymbol{\beta}}$, $\mathbf{E}_{\boldsymbol{\gamma}}$, $\mathbf{E}_{\boldsymbol{\theta}}$, $\mathbf{G}_{\boldsymbol{\alpha}}$, $\mathbf{G}_{\boldsymbol{\beta}}$) jointly in an end-to-end fashion, and we alternate optimization steps between the face model and the image discriminator. For the first 50,000 iterations, we minimize loss function (14) for all blocks in the face model, using a batch size of 16 and the Adam optimizer with learning rate 0.00001. We use a batch size of 16 and a learning rate of 0.00001 for the discriminator as well.

$$L_G + 1000 L_{\text{image}}^{\text{real}} + 10 L_{\text{identity}}^{\text{real}} + 100 L_{\text{landmark}}^{\text{real}}$$
$$+ 1.0 L_{\text{albedo}}^{\text{real}} + 10^{-5} L_{\text{lighting}}^{\text{real}} + 1.0 \|\boldsymbol{\alpha}\|_2^2 + 1.0 \|\boldsymbol{\beta}\|_2^2 \quad (14)$$

After training the network for 50,000 iterations, we fine tune $\mathbf{E}_{\boldsymbol{\beta}}$, $\mathbf{E}_{\boldsymbol{\gamma}}$, $\mathbf{G}_{\boldsymbol{\beta}}$ for another 50,000 iterations by freezing the weights of $\mathbf{E}_{\boldsymbol{\alpha}}$, $\mathbf{E}_{\boldsymbol{\theta}}$, $\mathbf{G}_{\boldsymbol{\alpha}}$ and discarding the landmark loss to further improve the reconstruction quality.

**Refinement.** Denoting the combined face-and-hair reconstruction as $\hat{\mathbf{x}}_{\mathbf{c}} := \hat{\mathbf{x}}_{\mathbf{f}} \odot \mathbf{M}_{\mathbf{f}} + \hat{\mathbf{x}}_{\mathbf{h}} \odot \mathbf{M}_{\mathbf{h}}$, the refined image as $\hat{\mathbf{x}} := \mathbf{R}(\hat{\mathbf{x}}_{\mathbf{c}})$, and the original face and hair as $\mathbf{x}' := \mathbf{x} \odot (\mathbf{M}_{\mathbf{f}} + \mathbf{M}_{\mathbf{h}})$, we employ the following loss function for the refiner:

$$L_G + 8.0 \|f_{\text{VGG}}(\mathbf{x}') - f_{\text{VGG}}(\hat{\mathbf{x}})\|_2^2$$
$$+ 10 \big(1 - \cos(f_{\text{id}}(\mathbf{x}'), f_{\text{id}}(\hat{\mathbf{x}}))\big) \quad (15)$$

where the second term stands for the VGG-16 perceptual loss (Simonyan and Zisserman 2014; Johnson, Alahi, and Fei-Fei 2016). With a batch size of 16, we alternately train the refiner and the discriminator for 500,000 iterations, with learning rate 0.0001 (using the Adam optimizer). To prevent overfitting, we randomly translate $\mathbf{x}'$ and $\hat{\mathbf{x}}_{\mathbf{c}}$ together (with horizontal and vertical translations uniformly sampled from the range $[-15, 15]$ pixels).

## Hair Manipulation Algorithm Details

**Warp field calculation.** In each iteration of our hair manipulation algorithm, we first identify the visible triangles of the given face mesh with its reference pose, and compute the center of each triangle by taking the average of its vertices. Then, we project these triangle centers onto the image plane under both the reference and the target pose. Using the correspondences between the two projections, we construct a 2D warp by calculating how much each of the projected triangle centers moves in pixel space as a result of the pose change. To complete the warp field, we use the technique described below.

For the vertical component of the warp field, we simply copy the vertical warp component from the nearest neighbor that was assigned a warp. For the horizontal component, we use a heuristic to assign a fixed horizontal warp to every pixel on the left edge of the image and a different fixed horizontal warp to every pixel on the right edge; then the horizontal component of the warp field for the entire image is simply interpolated from the assigned warps. In particular, when we rotate the faces clockwise (counter-clockwise) around the vertical axis, we extend a ray from the center of the 3D face to the left (right) perpendicular to the face's plane of symmetry and identify the 3D point on the ray whose projection lies on the leftmost (rightmost) edge of the image. Next, we calculate by how much this point's projection into the image plane moves when the head rotates, and we multiply this number by 3 to obtain the horizontal warp that we assign to the leftmost (rightmost) column of the warp field. For the rightmost (leftmost) column, we heuristically choose a displacement of 10 pixels to the right (left). Finally, we interpolate between the assigned pixels using linear interpolation to obtain the horizontal warp of every image pixel.

**Regularization of the hair mask.** After we obtain a complete warp field, we apply it to the reference hair mask and orientation. The orientation is regularized as part of the MichiGAN (Tan et al. 2020) pipeline, whereas we regularize the mask by projecting it onto a PCA basis that we calculate from a dataset of hair masks. In particular, we construct our hair mask dataset by randomly selecting 10,000 samples from our FFHQ training set and combining it with 10,000 masks that we obtain from the USC Hair Salon database (Hu et al. 2015). For the latter, we attach 3D hair models from the USC Hair Salon database to the FLAME mean head model, which we rotate around the vertical axis by an angle uniformly sampled from $[-90°, 90°]$. The hair masks are obtained by rendering these 3D shapes. Finally, after downsampling all masks to $64 \times 64$ resolution, we construct a PCA basis with 50 principal components, onto which we project the hair masks at each iteration of our reposing algorithm.

Table 2: Architecture of the shape generator $\mathbf{G}_{\boldsymbol{\alpha}}$. The output of the network is a UV-representation of 3D shape, where the three channels of the $256 \times 256$ output represent 3D offsets (in $x$, $y$, and $z$) from the FLAME mean head shape (Li et al. 2017).

| layer type | kernel size / stride | output shape | activation |
|---|---|---|---|
| linear | – | $1024 \times 1$ | none |
| reshape | – | $16 \times 8 \times 8$ | – |
| conv2d | $4 \times 4$ / 1 | $32 \times 8 \times 8$ | tanh |
| upsample | – | $32 \times 16 \times 16$ | – |
| conv2d | $4 \times 4$ / 1 | $64 \times 16 \times 16$ | tanh |
| upsample | – | $64 \times 32 \times 32$ | – |
| conv2d | $4 \times 4$ / 1 | $64 \times 32 \times 32$ | tanh |
| upsample | – | $64 \times 64 \times 64$ | – |
| conv2d | $4 \times 4$ / 1 | $64 \times 64 \times 64$ | tanh |
| upsample | – | $64 \times 128 \times 128$ | – |
| conv2d | $4 \times 4$ / 1 | $64 \times 128 \times 128$ | tanh |
| upsample | – | $64 \times 256 \times 256$ | – |
| conv2d | $4 \times 4$ / 1 | $3 \times 256 \times 256$ | tanh |

Table 3: Architecture of the refiner $\mathbf{R}$. We employ a U-Net (Ronneberger, Fischer, and Brox 2015) with skip connections between the encoder and decoder parts of the network. In all layers of the encoder, we use the LeakyReLU activation function with a negative slope of 0.2. All layers of the decoder use the ReLU activation function.

| layer type | kernel size / stride | output shape | activation |
|---|---|---|---|
| conv2d | $4 \times 4$ / 2 | $64 \times 128 \times 128$ | LeakyReLU |
| conv2d | $4 \times 4$ / 2 | $128 \times 64 \times 64$ | LeakyReLU |
| conv2d | $4 \times 4$ / 2 | $256 \times 32 \times 32$ | LeakyReLU |
| conv2d | $4 \times 4$ / 2 | $512 \times 16 \times 16$ | LeakyReLU |
| conv2d | $4 \times 4$ / 2 | $512 \times 8 \times 8$ | LeakyReLU |
| conv2d_transpose | $4 \times 4$ / 2 | $512 \times 16 \times 16$ | ReLU |
| conv2d_transpose | $4 \times 4$ / 2 | $256 \times 32 \times 32$ | ReLU |
| conv2d_transpose | $4 \times 4$ / 2 | $128 \times 64 \times 64$ | ReLU |
| conv2d_transpose | $4 \times 4$ / 2 | $64 \times 128 \times 128$ | ReLU |
| conv2d_transpose | $4 \times 4$ / 2 | $3 \times 256 \times 256$ | ReLU |

In this work, starting from the pose of the original face image, we invoke this algorithm sequentially by imposing a pose change of $5^\circ$ in each iteration, going all the way to the full profile ($\pm 90^\circ$) poses.

## Ablation Study

**Training and loss functions.** In our experiments, we observed that pretraining on synthetic data training is crucial for our method to work, since we observed stability issues when we started out by training on real data. For real data training, our experiments suggested that all loss functions except for equations (11)–(13) in the paper (all loss functions except for the albedo regularization, lighting regularization, and regularization of shape and albedo codes) are crucial for our method to achieve reasonable face reconstructions. When we omitted albedo and lighting regularizations, we observed that our method converges to a state in which the albedo reconstructions are washed out and the appearance of the face is mostly attributed to the lighting, which suggests that albedo and lighting regularizations are important to achieve better albedo and lighting disentanglement.

**Refinement.** In this section, we analyze the impact of the refiner on our reconstructions by providing qualitative and quantitative comparison of our results with vs. without the refiner. In addition, we compare our reconstructions with a state-of-the-art nonlinear 3D morphable model proposed by Tran and Liu (Tran and Liu 2019). We illustrate our qualitative comparisons in Figure 7. It is clear that our reconstructions are much more accurate and photorealistic

than those of Tran and Liu. We also observe a notable improvement in photorealism using our complete model (with refiner) vs. using our model without the refiner. To quantify our observations, we calculate a face recognition (FR) score as the average cosine similarity between the feature vectors extracted from the ArcFace face recognition network (Deng et al. 2019a) for the original and reconstructed images (from our test dataset). We also compute the structural similarity index measure (SSIM) and peak signal-to-noise ratio (PSNR) (Hore and Ziou 2010) between the original and reconstructed images. We quanitatively compare our model with vs. without the refiner in Table 4. In Table 5, we perform quantitative comparisons with Tran and Liu. To obtain the results in Table 5, we masked out the hair, background, clothing, and teeth for fair comparison with Tran and Liu (Tran and Liu 2019).

## Additional Experiments and Results

In this section, we provide additional experiments and qualitative results.

**Expression manipulation.** In Fig. 8, we present additional expression manipulation results, generated using the same method described in the main paper (see Fig. 4 in the main paper).

**Lighting manipulation.** In Fig. 9, we present additional lighting manipulation results, generated using the same method described in the main paper (see Fig. 4 in the main paper).

**Shape transfer.** In Fig. 10, we present additional shape transfer results, generated using the same method described in the main
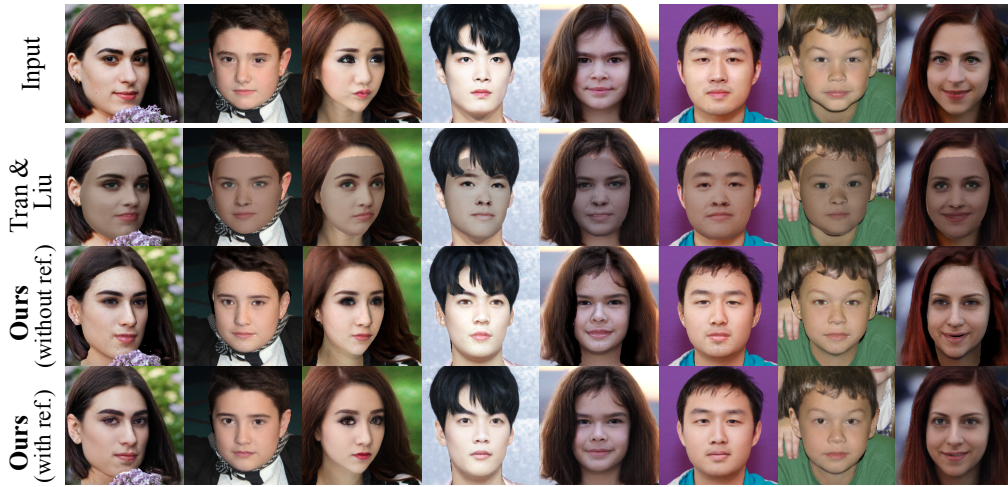
**Figure 7:** Ablation study and comparison with Tran and Liu (2019). On images from the test set (never seen during training), our method is able to reconstruct faces more accurately and photorealistically than Tran and Liu (2019). These results also demonstrate that our full model (with refinement) shows a notable improvement in image quality vs. our model without the refiner.

**Table 4:** Average face recognition (FR) scores, SSIM and PSNR between the original and reconstructed face images from our dataset. We observe a notable improvement due to the refinement.

|  | FR score ↑ | SSIM ↑ | PSNR ↑ |
| --- | --- | --- | --- |
| Ours (without refinement) | $0.66 \pm 0.10$ | $0.72 \pm 0.09$ | $22.24 \pm 2.63$ |
| **Ours (with refinement)** | $\mathbf{0.68 \pm 0.10}$ | $\mathbf{0.74 \pm 0.08}$ | $\mathbf{23.20 \pm 2.47}$ |

**Table 5:** Average face recognition (FR) scores, SSIM and PSNR between the original and reconstructed face images for our method (both with and without refinement) and Tran and Liu (2019). To obtain the results in this table, we masked out the hair, background, clothing, and teeth for fair comparison with Tran and Liu (2019). Our method achieves better scores in all three metrics.

|  | FR score ↑ | SSIM ↑ | PSNR ↑ |
| --- | --- | --- | --- |
| Tran and Liu (2019) | $0.51 \pm 0.12$ | $0.87 \pm 0.03$ | $20.96 \pm 1.57$ |
| Ours (without refinement) | $0.69 \pm 0.10$ | $0.87 \pm 0.04$ | $25.08 \pm 2.92$ |
| **Ours (with refinement)** | $\mathbf{0.71 \pm 0.10}$ | $\mathbf{0.88 \pm 0.04}$ | $\mathbf{26.17 \pm 2.71}$ |

paper (see Fig. 5 in the main paper), where we also compare with results obtained by a combination of StyleRig (Tewari et al. 2020a) and PIE (Tewari et al. 2020b). In this combination of previous methods, the image embedding is carried out by the optimization algorithm proposed in PIE, and the shape transfer is performed using StyleRig.

**Pose manipulation.** In Fig. 11, we present additional pose manipulation results, generated using the same method described in the main paper (see Fig. 6 in the main paper).

**Joint transfer of physical attributes.** Our face model's full disentanglement is demonstrated by its ability to transfer all physical attributes either individually or jointly. In Fig. 13, we present results of joint albedo and lighting transfer, as well as results of joint transfer of albedo, lighting, and shape.

**Interpolation in the latent space.** Although we do not impose any smoothness constraints within the latent spaces, the learned shape and albedo latent spaces enable smooth interpolation between different latent codes. We present our interpolation results in Fig. 12, where we simultaneously interpolate between the reconstructed shape code, albedo code, and lighting parameters of the reference and target images.

**Face anonymization.** Our face model can also be turned into a generative model for random faces by regularizing the latent code distributions during training. At each iteration, we calculate sample means and variances of the latent codes for shape and albedo over minibatches and regularize these statistics to match those of the standard Gaussian distribution by using a KL divergence loss. After this regularized training, we can randomly sample codes from a standard multivariate Gaussian distribution to produce realistic shapes and albedos. To enable random sampling of lighting conditions that match the distribution of lighting conditions present in the training set, we train a variational autoencoder on the reconstructed lighting parameters from the training set. In Fig. 14, we use random sampling of latent codes (i.e., random face generation) to anonymize face images from the test set. In each row of the figure, the input image is fed through our encoders to determine the pose and latent hair code. To anonymize the face, we randomly sample latent codes for shape, albedo, and lighting, while retaining the input image's background, pose, and hair code.

**Video of our results.** We also provide a video of our results where we demonstrate smooth manipulations in expression, lighting, albedo, shape, and their combinations.
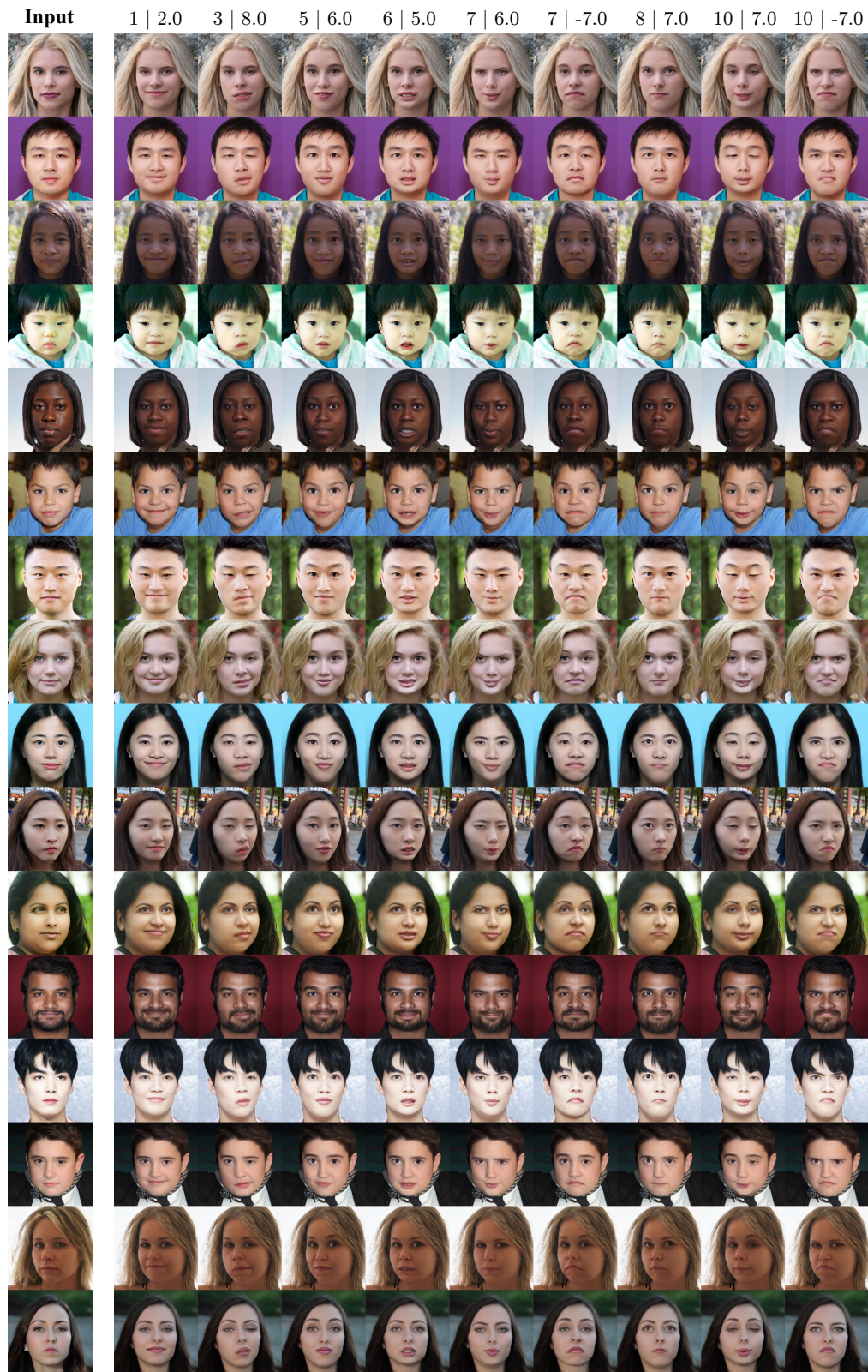
Figure 8: **Expression manipulation results.** We illustrate several expression changes with varying intensities. The two numbers above each column indicate which FLAME expression eigenvector is used and by how many standard deviations it is scaled.

Figure 9: **Lighting manipulation results.** For moderate variation, we rotate the reconstructed lighting around the camera axis by the angle listed above each column. For extreme lighting, we render the reconstructed 3D model using a point light source and Phong shading model.

Figure 10: **Shape transfer results and comparison**. We transfer the 3D shape of each source image to each target image while keeping everything else unchanged. Compared to the previous state of the art (StyleRig (Tewari et al. 2020a) + PIE (Tewari et al. 2020b)), our results (*top*) demonstrate more accurate shape transfer and much better disentanglement between shape and other face attributes (e.g., albedo, pose, and hair).
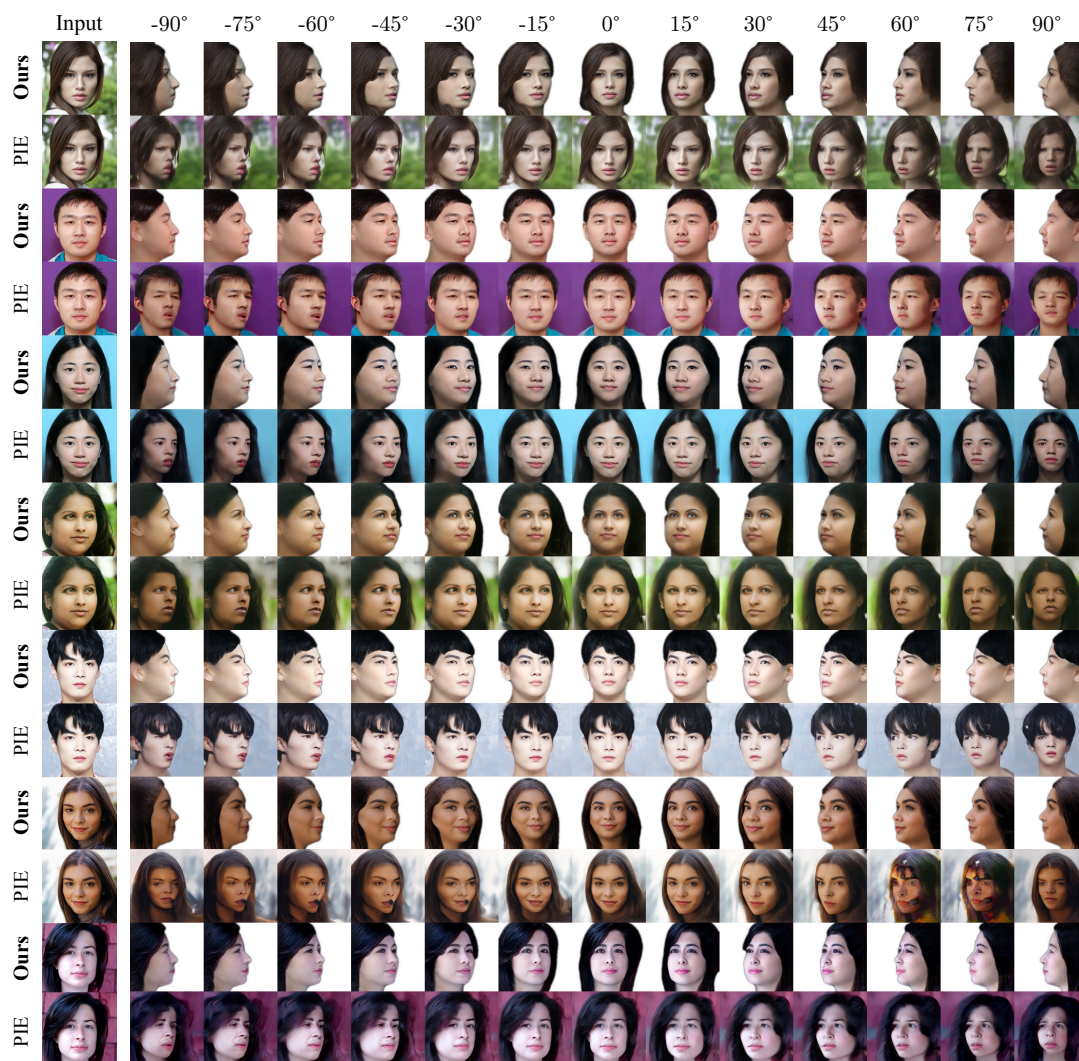
Figure 11: **Pose manipulation results and comparison with PIE** (Tewari et al. 2020b). To edit the pose of a given portrait image, we rotate the reconstructed faces in 3D and warp the hair in 2D. Our method is able to rotate portrait images all the way to profile pose while keeping the identity, expression, and illumination conditions unchanged.
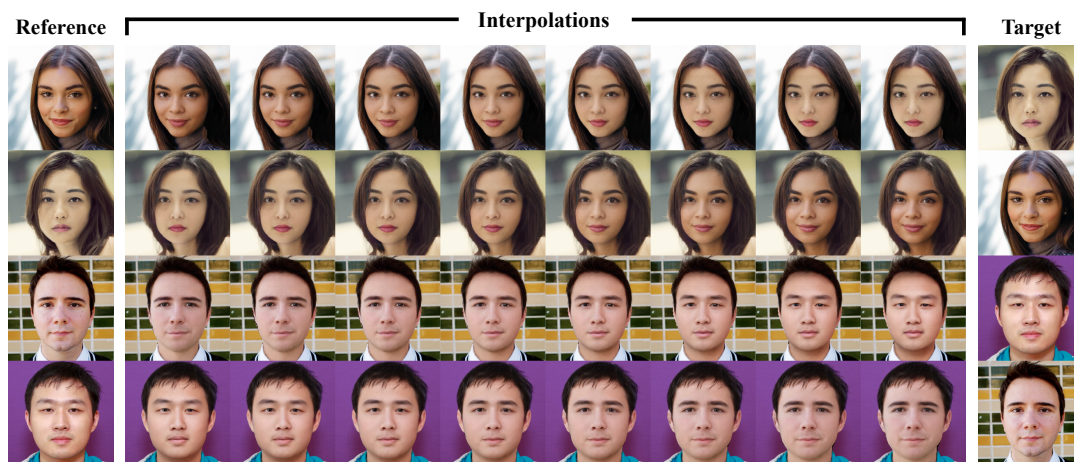


Figure 12: **Interpolation results.** Our method also allows for interpolation in the latent spaces. In each row, given a reference and a target image, we interpolate between their shape codes, albedo codes, and lighting parameters. For each interpolated image, the background and the latent code for hair are copied from the reference image.
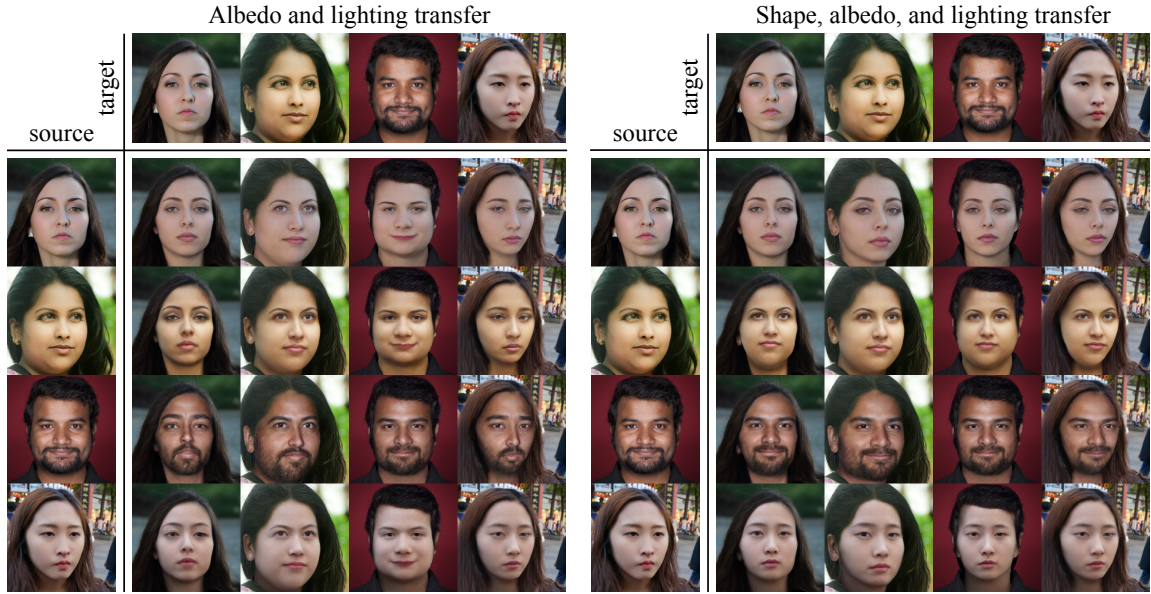
Figure 13: **Joint transfer of physical attributes.** Our method is able to transfer physical attributes jointly as well as individually. In this figure, we jointly transfer the indicated physical attributes of each source image to each target image while keeping the other parameters of the target image unchanged. *Left:* Albedo and lighting transfer. *Right:* Shape, albedo, and lighting transfer.



Figure 14: **Face anonymization results.** Our model can also be used to sample novel faces by regularizing the latent code distributions during training, which can be used for face anonymization.