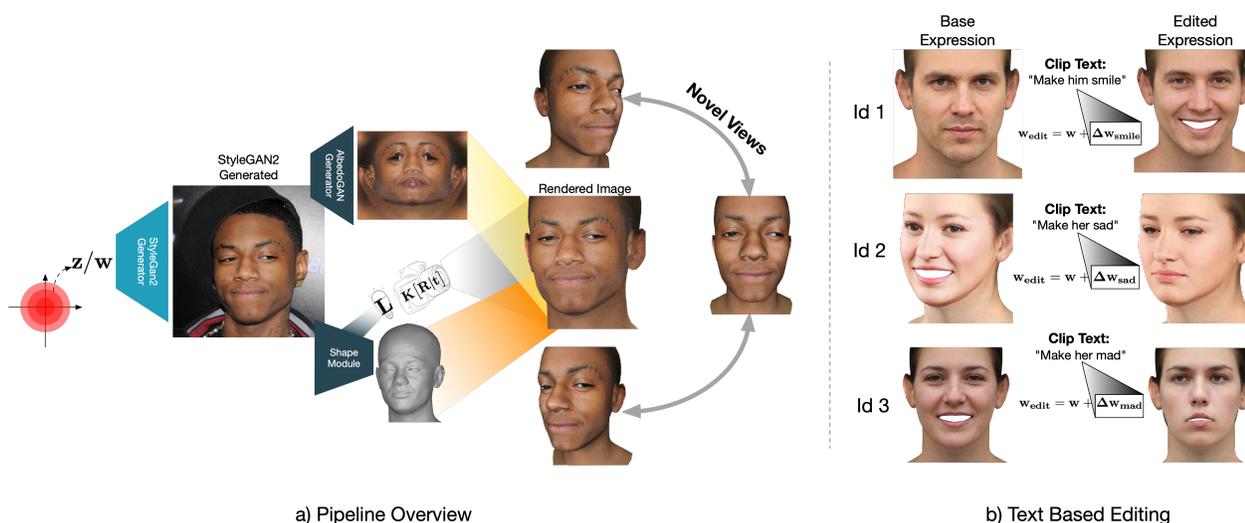


# Towards Realistic Generative 3D Face Models

Aashish Rai<sup>1</sup>    Hires Gupta<sup>\*1</sup>    Ayush Pandey<sup>\*1</sup>    Francisco Vicente Carrasco<sup>1</sup>  
 Shingo Jason Takagi<sup>2</sup>    Amaury Aubel<sup>2</sup>    Daeil Kim<sup>2</sup>    Aayush Prakash<sup>2</sup>  
 Fernando De la Torre<sup>1</sup>

<sup>1</sup>Carnegie Mellon University    <sup>2</sup>Meta Reality Labs

<https://aashishrai3799.github.io/Towards-Realistic-Generative-3D-Face-Models>



a) Pipeline Overview

b) Text Based Editing

Figure 1: 3D generative face model. a) High-resolution 3D shape and albedo recovered from a StyleGAN2 generated image. Novel views can be rendered using the estimated face model. b) Editing of 3D faces with text. This method allows for 3D expression manipulation through guidance with the CLIP model.

## Abstract

In recent years, there has been significant progress in 2D generative face models fueled by applications such as animation, synthetic data generation, and digital avatars. However, due to the absence of 3D information, these 2D models often struggle to accurately disentangle facial attributes like pose, expression, and illumination, limiting their editing capabilities. To address this limitation, this paper proposes a 3D controllable generative face model to produce high-quality albedo and precise 3D shape leveraging existing 2D generative models. By combining 2D face generative models with semantic face manipulation, this method enables editing of detailed 3D rendered faces. The proposed framework utilizes an alternating descent optimization approach over shape and albedo. Differentiable rendering is used to train high-quality shapes and albedo

without 3D supervision. Moreover, this approach outperforms most state-of-the-art (SOTA) methods in the well-known NoW and REALY benchmarks for 3D face shape reconstruction. It also outperforms the SOTA reconstruction models in recovering rendered faces' identities across novel poses by an average of 10%. Additionally, the paper demonstrates direct control of expressions in 3D faces by exploiting latent space leading to text-based editing of 3D faces.

## 1. Introduction

The success of language models like GPT-3 [41], and more recently, the release of text-to-image models like GLIDE [35], DALLE-2 [43], or Imagen [45] have all contributed to the enormous popularity of generative AI. Besides generating images with unprecedented visual quality, these models also show remarkable generalization ability to novel texts with complex compositions of concepts, mak-

\* equal contribution

ing them generalists for image synthesis. In the context of faces, StyleGAN2 [28] has the capability of generating powerful face images that are frequently indistinguishable from reality. While these 2D generative models create high-quality faces for many applications of interest, such as facial animation [27, 55], expression transfer [30, 52, 36] virtual avatars [34], these 2D models often encounter difficulties when it comes to effectively disentangle facial attributes like pose, expression, and illumination. As a result, their capacity to edit such attributes is limited. Moreover, a 3D representation (shape, texture) is crucial to many entertainment industries—including games, animation, and visual effects—that are demanding 3D content at increasingly enormous scales to create immersive virtual worlds. Recall that many applications of interest require 3D assets that are consumable by a graphics engine (e.g., Unity [54], Unreal [11]).

To address this demand, recently, researchers have proposed generative models to generate 3D faces [1, 50, 19]. Even though these algorithms perform well, the lack of diverse and high-quality 3D training data has limited the generalization of these algorithms and their use in real-world applications [53]. Another line of research involves using parametric models like 3DMM[3], BFM[39], FLAME[3], and derived methods [15, 32, 10, 51, 46] to approximate the 3D geometry and texture of a 2D face image. While these 3D face reconstruction techniques can reasonably recover low-frequency details, they typically do not recover high-frequency details. Also, predicting high-resolution texture maps that capture details remains an unaddressed problem. Most of the works focusing in this direction either emphasize mesh or texture. However, a generative 3D face model that can generate both high-quality texture and a detailed mesh with the same quality as 2D models is still missing.

This paper proposes a 3D generative model for faces using a self-supervised approach that can generate high-resolution texture and capture high-frequency details in the geometry. The method leverages a pretrained StyleGAN2 to generate high-quality 2D faces (see Fig. 1 a). We propose a network, AlbedoGAN, that generates light-independent albedo directly from StyleGAN2’s latent space. For the shape component, the FLAME model [33] is combined with per-vertex displacement maps guided by StyleGAN’s latent space, resulting in a higher-resolution mesh. The two networks for albedo and shape are trained in alternating descent fashion. The proposed method outperforms SOTA methods in shape estimation, such as DECA [15] and MICA [61], by 20% and 1.1%, respectively. It’s worth noting that MICA only generates a neutral and frontal smooth mesh, while the proposed algorithm can generate any expression. Fig. 1(a) shows how an image generated by StyleGAN2 can be uplifted to 3D with a detailed shape and albedo, being able to render realistic 3D faces. Finally,

given the 3D face asset, our algorithm can edit the face in 3D. For example, Fig. 1(b) illustrates expression manipulation through text-based editing guided by the CLIP model [40]. Briefly stated, our main contributions are:

1. A self-supervised method to leverage StyleGAN2 into a 3D generative model producing high-quality albedo and a detailed mesh. We introduce AlbedoGAN, a single-pass albedo prediction network that generates high-resolution albedo and decouples illumination using shading maps.
2. We show that our model outperforms existing methods in capturing high-frequency facial details in a mesh. Moreover, the proposed method reconstructs 3D faces that recover identity better than SOTA methods.
3. We propose a displacement map generator capable of decoding per-vertex displacements directly from StyleGAN’s latent space using detailed normals of the mesh.
4. Since our entire architecture can generate 3D faces from StyleGAN2’s latent space, we can perform face editing directly in the 3D domain using the latent codes or text.

## 2. Related Work

Reconstructing a 3D Face from a single 2D image is an ill-posed problem that has intrigued researchers for decades. Blanz et al. [3] took the first significant step in this direction when they introduced 3D Morphable Models (3DMM) [3] in 1999 as general face representation. Their work inspired decades of work in estimating parameters for 3DMM to find a textured mesh that best fits an input 2D face. While estimating parameters for parametric models like 3DMM and its advanced versions like FLAME [33] has been the bulk of the focus for researchers, there has also been a good amount of work done in learning volumetric representations (e.g., NeRF) for a face. In particular, the focus has been on using Neural Implicit Representation [24, 17, 59, 44, 5, 60] to learn density and radiance to represent a face. However, it has been widely noticed that they are prone to generating artifacts and consume a lot of time to render these detailed representations. Furthermore, these methods do not generate a topologically uniform mesh, and therefore do not directly serve applications in graphic engines, face animation, avatar creation, etc. Due to the above-mentioned reasons, we do not consider implicit representation in our research. In the following sections, we describe work that is related to individual components of our framework.

### 2.1. Texture Generation for 3D Mesh

Most of the work done in synthesizing texture for a mesh can be broadly divided into two parts: using a parametric texture model like 3DMM or Basel Face Model (BFM) [15, 10, 32, 51] or a GAN-based approach [18, 20, 21, 23, 56, 31, 22, 7] to generate texture. Using a parametric model BFM [39], works like [15, 10, 32, 51] learn an encoder to predict the parameters that generate a texture that best fits

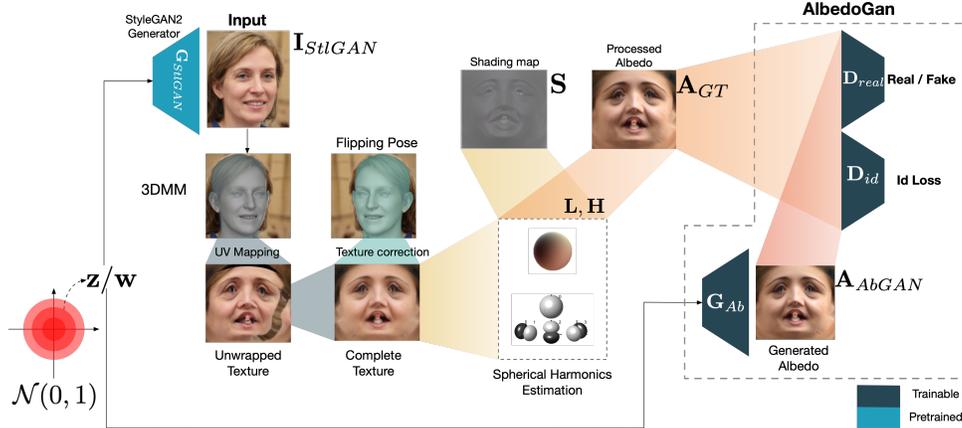


Figure 2: AlbedoGAN. Pose-invariant albedo,  $\mathbf{A}_{GT}$ , obtained by texture extraction and synthesis 3.1, is used to train StyleGAN2 generator,  $\mathbf{G}_{Ab}$ , for a given latent code  $\mathbf{z}/\mathbf{w}$ . We use a 3DMM fitting, image blending, and Spherical Harmonics to extract  $\mathbf{A}_{GT}$ .  $\mathbf{D}_{real}$  and  $\mathbf{D}_{id}$  are introduced to generate realistic images and identity consistent albedo, respectively.

the visible part of the face. Since the model uses representation in low dimensional PCA-based space, they generate an approximate texture that often lacks high-frequency details, which correspond to low variation direction in the projected latent space and do not lead to photo-realistic rendering. Recently, with the advent of GANs, there are works that leverage them to extract texture [18, 20, 31]. TBGAN [20] trains a Progressive GAN to generate a high-quality texture with a differentiable renderer in a supervised setting. On the other hand, OSTeC [18] uses 3DMM as initialization to generate a textured mesh and then uses GAN inversion with StyleGAN2 to generate multiple views of this mesh and extract texture. This is an extremely time-consuming and takes several minutes per image. We propose AlbedoGAN, which is able to generate high-quality texture in a single pass in a time-efficient manner. AlbedoGAN generates textures that maintain identity over multiple poses - where most previous methods struggle.

## 2.2. 3D Shape prediction from 2D Image

Similar to texture, Blanz et al.’s seminal work [3] can represent a face mesh in a low-dimensional PCA-based space. This led to the development of a huge corpus of work [10, 51, 46, 32] in 3D face reconstruction focused on learning an encoder that could predict parameters for generating shapes using 3DMM, given a 2D image. The encoder could be learned in a self-supervised way with 2D image losses [15, 32, 10, 51, 46], or with 3D supervision, [61]. After 3DMM[3], there have been newer parametric models BFM [39] and FLAME [33], which have been learned from more subjects and encode the structure of the face better. The explosive development of 3D Face Reconstruction led to the creation of NoW benchmark [46] as a common ground for comparison across 3D Face reconstruction methods. Currently, DECA [15], and MICA [61] show best reconstruc-

tion on NoW Benchmark. DECA’s architecture is inspired by RingNet [46]. However, instead of 3DMM, DECA uses FLAME [33], and it adds an encoder-decoder network that learns to generate displacement maps to learn animatable details in UV space. MICA [61], the current state-of-art model, leverages ArcFace backbone [8] to learn the actual shape of the face by regressing it on a high-quality 3D face scan dataset[6, 16, 4].

## 3. AlbedoGAN Training

Albedo constitutes one of the crucial parts of a 3D face model, since face appearance is largely dictated by it. To generate high quality 3d models, we need to generate albedo that generalize over pose, age, and ethnicity. However, training such a diverse albedo generative model requires a massive database of 3D scans, which is neither cost nor time effective. An efficient way of extracting textures from existing 2D images is fitting a 3DMM and capturing a UV mapped texture. However, this ”pseudo” texture does not generalize well over poses nor disentangle shadows. In this paper, we leverage 3DMM fitting, image blending, and Spherical Harmonics lighting to capture high-quality  $1024 \times 1024$  resolution albedo that generalizes well over different poses and tackles shading variations.

This section describes albedo training, refer Fig. 2 for an overview. The first step includes texture extraction and correction, 3.1, followed by the use of a spherical harmonics model to extract albedo from texture (Section 3.2). Section 3.3 explains training AlbedoGAN - a StyleGAN2 model to generate albedo corresponding to the given latent code  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ .

### 3.1. Texture Extraction and Correction

First, we establish a correspondence between the input 2D image  $\mathbf{I} \in \mathbb{R}^{w \times h \times 3}$  and the UV domain by taking orthogonal projection of a mesh fitted on the given image using a 3DMM [33]. Using this correspondence, we get the RGB values for the UV texture and perform barycentric interpolation to fill out the missing pixels.

Now, texture correction is performed to fill the occluded areas by leveraging pose information. This happens by projecting the flipped input image and fitted mesh, and collecting the pixels corresponding to the missing parts. These pixels are blended to the original texture to get a complete pose-invariant texture.

### 3.2. Albedo from Texture

As shown in Fig. 2, the next step includes obtaining an albedo and shading map from the unevenly illuminated texture. Following previous works [15, 10], we made the following assumptions: (1) The illumination model is Spherical Harmonics (SH), (2) light source is distant and monochromatic, and (3) surface reflectance is Lambertian. The shading map can thus be calculated as

$$\mathbf{S}_{ij} = \sum_{b=1}^9 \mathbf{L}_b \mathbf{H}_b(\mathbf{N}_{ij})$$

where,  $\mathbf{H}_b : \mathbb{R}^3 \rightarrow \mathbb{R}$  are the SH basis function,  $\mathbf{L}_b \in \mathbb{R}^3$  are SH coefficients, and  $\mathbf{N}_{ij} \in \mathbb{R}^3$  are surface normals. The relation between albedo  $\mathbf{A} \in \mathbb{R}^{w \times h \times 3}$ , texture  $\mathbf{T} \in \mathbb{R}^{w \times h \times 3}$ , and shading map  $\mathbf{S} \in \mathbb{R}^2$  can then be defined as  $\mathbf{T}_{ij} = \mathbf{A}_{ij} \odot \mathbf{S}_{ij}$ , where,  $\odot$  is the Schur product.

### 3.3. Training

In this section, we address the AlbedoGAN model training procedure. Later, the resulting model will be fine-tuned during the shape and displacement map training process, taking into account geometry and more sophisticated Phong illumination model.

Our approach, AlbedoGAN, is built upon a generative model that can synthesize face images corresponding to a latent vector, to this end we selected StyleGAN2. This model is best suited to our requirements as it uses a mapping network  $M_L$  that maps an input noise vector  $\mathbf{z} \in \mathbb{R}^{512}$  to an intermediate latent vector  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ . This mapping,  $\mathbf{w} = M_L(\mathbf{z})$ , adds the ability for manipulation and better projection. Consequently, we use  $\mathbf{w}$  as latent space for AlbedoGAN.

Hence, face images are generated by randomly sampling  $\mathbf{w}$  using a pretrained StyleGAN2. The generated images act as input,  $\mathbf{I}_{StuGAN}$ , see Fig.2. The same latent codes,  $\mathbf{w}$ , are used in the AlbedoGAN generator, to produce  $\mathbf{A}_{AbGAN}$ . As shown in Fig.2,  $\mathbf{I}_{StuGAN}$  is passed through the texture extraction, 3.1 and albedo extraction, 3.2, steps to obtain  $\mathbf{A}_{GT}$ .

This albedo is used as a ground truth for the training of the real/fake discriminator,  $\mathbf{D}_{real}$ . Additionally, we constraint AlbedoGAN to generate identity consistent albedos by introducing an identity discriminator,  $\mathbf{D}_{id}$ . To this intent, we use the features of a pretrained face recognition model [26]  $F : \mathbb{R}^{w' \times h' \times 3} \rightarrow \mathbb{R}^{512}$ . Our identity loss is defined as cosine distance between the identity features of predicted albedo  $\mathbf{A}_{AbGAN}$  and  $\mathbf{A}_{GT}$  as:

$$L_{id}(\mathbf{A}_{AbGAN}, \mathbf{A}_{GT}) = 1 - \frac{F(\mathbf{A}_{AbGAN}) \cdot F(\mathbf{A}_{GT})}{\|F(\mathbf{A}_{AbGAN})\|_2 \|F(\mathbf{A}_{GT})\|_2} \quad (1)$$

## 4. Alternating Descent in Albedo and Shape

In this section, we describe our regression method to 3D shape given a face image and the Differentiable Rendering, DR, based approach to fine-tune AlbedoGAN. This fine-tuning process for AlbedoGAN takes into account expression, camera pose, and the Phong illumination model.

Unfortunately, jointly optimizing all the components (shape, albedo, illumination, etc.) that produce the best rendered face that is consistent with the input image is computationally expensive. Thus, we propose using Alternating Descent for optimization. First, we optimize the shape for a few iterations while freezing AlbedoGAN. Next, AlbedoGAN is fine-tuned using the updated 3D shape, with more detailed normals of the shape. This alternating optimization cycle is repeated throughout the course of the training process. Next, we first describe albedo optimization, 4.1, followed by shape optimization, 4.2.

### 4.1. Albedo optimization

To fine-tune AlbedoGAN using the information of the 3D shape and illumination model, we first assume we have an estimate of a detailed 3D shape  $\mathbf{M}'$ . As shown in Fig. 3, given an estimated mesh  $\mathbf{M}'$ , predicted albedo  $\mathbf{A}_{AbGAN}$ , pose  $\mathbf{p}$ , and light  $\mathbf{l}$ , we can generate a detail rendered image  $\mathbf{I}_{ren}$  using DR,  $R$  as:

$$\mathbf{I}_{ren} = R(\mathbf{M}', \mathbf{A}_{AbGAN}, \mathbf{p}, \mathbf{l})$$

The overall loss function  $\mathcal{L}$  is defined as a sum of the following terms:

$$\mathcal{L} = \lambda_{sym-rec} \mathcal{L}_{sym-rec} + \lambda_{id} \mathcal{L}_{id} + \lambda_{perc} \mathcal{L}_{perc} + \lambda_{lmk} \mathcal{L}_{lmk}$$

Where each loss is defined as follows:

**Symmetric Reconstruction Loss,  $\mathcal{L}_{sym-rec}$ :** A simple supervision function that encourages low-level similarity in the predicted image and the corresponding ground truth and symmetry in the estimated albedo. We use the Mean

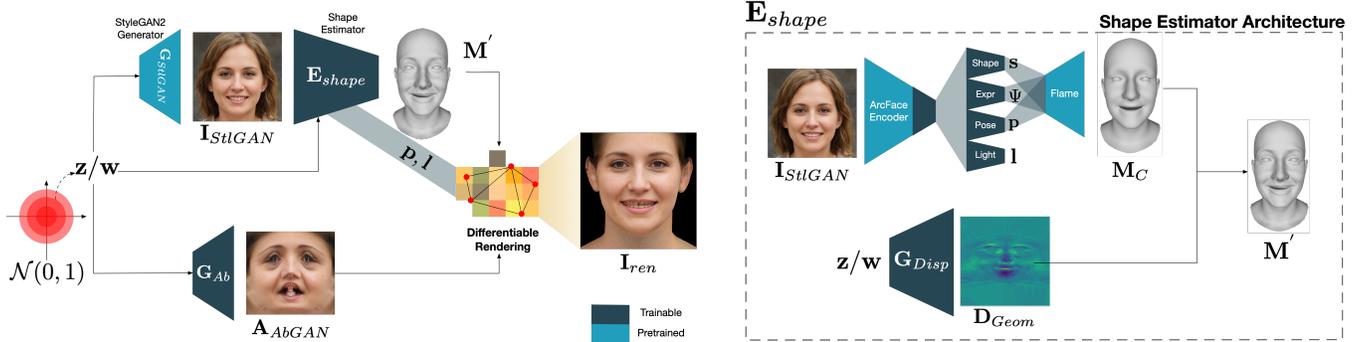


Figure 3: Overview of our generative model. The AlbedoGAN generator,  $G_{Ab}$ , is used to synthesize albedo  $A_{AbGAN}$  corresponding to a latent code  $w$ .  $G_{StlGAN}$  generates a 2D image,  $I_{StlGAN}$ , given to the shape model,  $E_{shape}$ , to get a detailed mesh,  $M'$ . Finally, a differentiable renderer (DR) is used to synthesize  $I_{ren}$  from the mesh  $M'$ , albedo  $A_{AbGAN}$ , lighting  $l$ , and pose  $p$ . Losses between  $I_{ren}$  and  $I_{StlGAN}$  are used to train the shape model and the AlbedoGAN via Alternating Descent.

Squared Error (MSE) to calculate reconstruction error.

$$\mathcal{L}_{sym\_rec}(I_{ren}, I_{StlGAN}) = \text{MSE}(I_{StlGAN}, I_{ren}) + \underbrace{\text{MSE}(I'_{StlGAN}, I'_{ren})}_{\text{Albedo symmetry consistency term}}$$

Where  $I'_{StlGAN}$  is the flipped ground truth obtained through StyleGAN2, and  $I'_{ren}$  is the rendered estimated image through flipping  $A_{AbGAN}$ , pose  $p$ , and light  $l$ .

**Identity Loss,  $\mathcal{L}_{id}(I_{ren}, I_{StlGAN})$ :** This loss term is introduced with the intent of making the AlbedoGAN generator learn to match the identity of the rendered face,  $I_{ren}$ , with the ground truth,  $I_{StlGAN}$ . As in section 3.3, we use a pretrained face recognition model [26] for feature extraction. The cosine distance, eq. 1, is used to calculate the identity loss while fine-tuning the model.

**Perceptual Loss,  $\mathcal{L}_{perc}$ :** This perceptual loss is introduced with the goal of forcing AlbedoGAN to generate a  $A_{AbGAN}$  that matches the visual appearance of  $I_{StlGAN}$ . Motivated by the existing research, we selected a VGG16 based feature extractor [47], a pretrained face recognition model. We use the output of *relu3\_3* as the image features. The loss is calculated by the L2 distance between the feature vectors from  $I_{ren}$  and  $I_{StlGAN}$ .

**Landmark Loss,  $\mathcal{L}_{lmk}$ :** AlbedoGAN is also fine tuned using 68 facial landmarks detected on the ground truth and the rendered image to avoid misaligned generations. We used a SOTA face landmark detection [57] to predict 68 landmarks on  $I_{StlGAN}$  and  $I_{ren}$ . The loss is calculated using MSE between the two set of landmarks.

## 4.2. Shape Model and Optimization

We proceed with the shape optimization, while freezing AlbedoGAN. We sample a latent vector  $w \in \mathbb{R}^{18 \times 512}$  and use a pretrained StyleGAN2 model,  $G_{StlGAN}$ , to generate a 2D face image,  $I_{StlGAN}$ , and AlbedoGAN generator,

$G_{AbGAN}$ , to generate albedo,  $A_{AbGAN}$ . Figure 3 describes the detailed architecture of our shape model,  $E_{shape}$ . We leverage ArcFace backbone [9] to predict the face shape ( $s$ ), expression ( $\psi$ ), lighting ( $l$ ), and camera pose ( $p$ ) parameters for the given image  $I_{StlGAN}$ . The shape embedding vector  $s \in \mathbb{R}^{300}$ , pose  $p \in \mathbb{R}^6$ , and expression  $\psi \in \mathbb{R}^{50}$  parameters are fed into a parametric face model that gives us a coarse mesh representation ( $M_c$ ) as described below:

$$M_c(s, \psi) = T + B_s s + B_\psi \psi \quad (2)$$

where  $M_c$  represents the generated coarse mesh synthesized by a 3DMM decoder. Specifically, we use FLAME [33] as our mesh decoder, which generates a coarse mesh with  $N = 5023$  vertices. This coarse mesh is computed by using a template mesh,  $T \in \mathbb{R}^{3N}$ , representing a mean human face and different principal components  $B_s \in \mathbb{R}^{3N \times 300}$  and  $B_\psi \in \mathbb{R}^{3N \times 50}$  corresponding to the shape and expression terms respectively.

To capture high-frequency details in meshes, we learn a displacement generator  $G_{Disp}$  to augment the coarse mesh,  $M_c$ , with a detailed UV displacement map  $D_{Geom} \in [-0.01, 0.01]^{n \times n}$ . Recent research [2, 58] has shown that StyleGAN's latent space contains information about high-frequency details of a face. Using this insight, we predict displacement maps to capture expression and pose dependent per vertex offsets. The latent code  $w$  is the same as used in AlbedoGAN and the StyleGAN2 model. Finally, we combine the displacement map along the vertex normals of the mesh  $M_c$  to get a detailed mesh  $M'$  by adding them in the UV domain.

We use the detailed mesh,  $M'$ , along with the predicted pose  $p$ , light  $l$  parameters, and the synthesized albedo  $A_{AbGAN}$  to render an image  $I_{ren}$  as described below:

$$I_{ren} = R(M_c, A_{AbGAN}, p, l) \quad (3)$$

We apply multiple 2D image-based losses, including identity loss, perceptual loss, and landmark loss between  $\mathbf{I}_{StlGAN}$  and  $\mathbf{I}_{ren}$  to optimize the mesh representation in a self-supervised fashion. In addition to the losses, we also calculate a shape center loss eq. 4 on images belonging to identity  $i$ . In particular, eq. 4 tries to reduce the distance between shape vector for all the images and their corresponding mean  $\mu_i$ . Besides this, we also perform L2 regularization, eq. 5, on predicted shape  $\mathbf{s}$ , expression  $\psi$ , and displacement maps  $\mathbf{D}_{Geom}$  that enforce a prior distribution towards the mean face.

$$\mathcal{L}_{sc} = \sum_{i=0}^N \sum_{k=0}^K \|\mathbf{s}_{i,k} - \mu_i\|_2^2 \quad (4)$$

$$\mathcal{L}_{reg} = \|\mathbf{s}\|_2^2 + \|\psi\|_2^2 + \|\mathbf{D}_{Geom}\|_F^2 \quad (5)$$

The overall loss function  $\mathcal{L}$  is defined as a weighted sum:

$$\mathcal{L} = \lambda_{id}\mathcal{L}_{id} + \lambda_{perc}\mathcal{L}_{perc} + \lambda_{lmk}\mathcal{L}_{lmk} + \lambda_{sc}\mathcal{L}_{sc} + \lambda_{reg}\mathcal{L}_{reg}$$

## 5. Experiments

This section describes the implementation details, quantitative, and qualitative evaluation of the shape and texture reconstruction models, along with SOTA comparison.

### 5.1. Dataset

We randomly sampled 100K, 512-dimensional random vectors  $\mathbf{z} \in \mathbb{R}^{512}$  from a Gaussian distribution and generated the corresponding  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$  from the StyleGANs mapping network as  $\mathbf{w} = M_L(\mathbf{z})$ . These intermediate latent vectors  $\mathbf{w}$  are used to generate 100K images  $\in \mathbb{R}^{1024 \times 1024 \times 3}$  from a pretrained StyleGAN2 [29] generator. To ensure diversity in the generated images across ethnicity, expression, age, and pose; we followed the work in [42]. The texture-preprocessing step (as described in Sec. 3.1) is used to get the complete GT-albedo corresponding to all the samples in the dataset.

To train our shape model with 2D images, we chose 30K of the previously sampled  $\mathbf{z}$  vectors. Then we perform latent space editing in the  $\mathbf{w}$  space to generate a total of 11 images (belonging to different expressions and poses) per identity using StyleGAN2 implementation [42] of InterFaceGAN [48]. We estimate 68 landmarks on all the GT images using the FAN [57] landmark detection.

### 5.2. Implementation Details

**Albedo Generation:** PyTorch [37] is used as the implementation framework on CUDA enabled system with NVIDIA RTX A4500 GPUs. We use the PyTorch implementation of StyleGAN2<sup>1</sup> for albedo and image generation. To generate face images from StyleGAN2, we use the

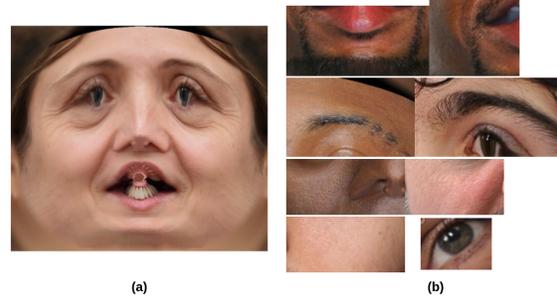


Figure 4: (a) Randomly generated albedo from AlbedoGAN. (b) Patches of randomly generated albedo. AlbedoGAN can generate high-quality albedo of 1K resolution.

official pretrained weights trained on the FFHQ dataset for 1024 × 1024 resolution. We use Adam optimizer to train the AlbedoGAN with learning rate  $\alpha_{gen} = \alpha_{disc} = 2e^{-3}$  and  $\beta_1 = 0, \beta_2 = 0.99$ . The generator was regularized after every 4 iteration, while the discriminator after every 16 training iterations.

In the first step, we train the AlbedoGAN from GT-albedo and use the GAN loss and ID loss eq.1 for supervision. Similar to StyleGAN2, for the GAN loss, we use the element-wise Softplus, which can be defined as  $Softplus(x) = \frac{1}{\beta} * \log(1 + \exp(\beta * x))$ . The  $\mathcal{L}_{ID}$  is calculated between the predicted albedo from the generator  $\mathbf{A}_{AbGAN}$  and the GT-albedo  $\mathbf{A}_{GT}$ . The  $\lambda_{ID}$  was set to 1 during this training. We trained the model on 8 GPUs, and it took around 32 hours for complete training (batch size 32).

This gave us a robust albedo generator capable of generating a pose-invariant albedo corresponding to a given  $\mathbf{z} \in \mathbb{R}^{512}$  or  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ .

**Optimizing Shape:** We train our shape model,  $\mathbf{E}_{shape}$ , on the synthetically generated images by StyleGAN2 capturing multiple images of the same identity across varying expression & pose. We run a face detector [25] on the input images and scale the face crops to a resolution of 224 × 224 before passing them to our shape model. The shape model consist of an ArcFace backbone that is initialized to the weights learned by [61] and a convolution-styled decoder ( $\mathbf{G}_{Disp}$ ) respectively. The whole pipeline is optimized using Adam Optimizer with a learning rate of  $1e^{-4}$ . The final loss is calculated between rendered images  $\mathbf{I}_{ren}$  and GT  $\mathbf{I}_{StlGAN}$ , where  $\lambda_{id}$  is set to 0.5, and  $\lambda_{perc}, \lambda_{lmk}, \lambda_{sc}$  and  $\lambda_{reg}$  are set to 1, 5, 1 and  $1e^{-4}$  respectively.

**Fine-tuning AlbedoGAN using DR:** Once we have a pretrained AlbedoGAN and a good shape estimator, we now fine-tune the AlbedoGAN. This makes the albedo generator learn to capture more details from the GT-face. PyTorch3D is used as the differentiable renderer in all our experiments under this section. We kept using  $\mathcal{L}_{gan}$  from previous AlbedoGAN training but gave more gravity to the

<sup>1</sup><https://github.com/rosinality/stylegan2-pytorch>

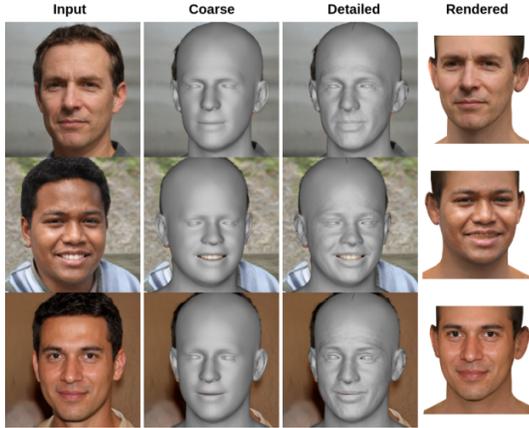


Figure 5: Randomly generated coarse mesh, detailed mesh and rendered faces from our model, for input 2D faces.

rendering losses. The rendering loss is calculated between  $\mathbf{I}_{StlGAN}$  and  $\mathbf{I}_{ren}$  along with a symmetric reconstruction loss between  $\mathbf{A}_{GT}$  and  $\mathbf{A}_{AbGAN}$  to maintain low-level consistency in the fine-tuned albedo.  $\lambda_{sym.rec}$  and  $\lambda_{ID}$  were set to 0.1 for albedo based losses.  $\lambda_{sym.rec}$ ,  $\lambda_{ID}$ ,  $\lambda_{perc}$  and  $\lambda_{lmk}$  were set to 1, 1, 1, 0.5 for rendering losses, respectively. We fine-tuned the model for another 48 hours, keeping the batch size of 8. The results of randomly generated textures from our AlbedoGAN are shown in Fig. 4.

Once we fine-tune the AlbedoGAN, we again train the shape model with updated albedo weights and repeat this step multiple times until we get a final model that can synthesize realistic-looking 3D faces corresponding to the given 2D images.

Fig. 1(a), 5, 6 shows the reconstructed 3D faces generated from our model for multiple poses. It is interesting to see how our model generalizes well over different poses and generates realistic-looking 3D faces. Section B.3 in supplementary demonstrates lighting and pose control in rendered faces using corresponding parameters and shading maps. Some additional results on testing our pipeline on real-world images using GAN inversion can be seen in supplementary section C.

### 5.3. Evaluation of Shape and Texture

#### 5.3.1 NoW Benchmark - shape reconstruction

NoW benchmark [46] is a standard benchmark to evaluate the accuracy of 3D meshes estimated from 2D images. It consists of 2054 images for 100 test subjects across different expressions, poses, and occlusions, split across two sets for validation (20 subjects) and test (80 subjects). NoW provides 3D ground truth meshes for each test subject, and the predicted mesh is rigidly aligned with the ground truth mesh using 3D face landmarks. The per-vertex error is then used for all the subjects to compute the mean, median, and standard deviation of the errors. Table 1 depicts the com-

parison of our model with the current SOTA methods, including DECA and MICA. Fig. 6 illustrates the visual comparison among these methods. Our model outperforms the DECA model by achieving a **23%** better median error in coarse mesh and a **20%** better median error in the detailed mesh. Our approach can reconstruct realistic-looking rendered faces, and model accurate head shapes, especially for faces with big heads. We also observe an improvement over the MICA model that was trained on 3D face scans [6, 16, 4] with 2300 subjects on the NoW validation set. As illustrated in 6, our method produces a more detailed mesh, capturing wrinkles, expression, pose and head shape correctly by only training on synthetic images.

Table 1: Reconstruction error on the NoW Benchmark.

Method	Median (mm)	Mean (mm)	Std (mm)
<b>Validation Set</b>			
Deep3D [10]	1.286	1.864	2.361
DECA [15]	1.178	1.464	1.253
MICA [61]	0.913	1.130	<b>0.948</b>
Ours	<b>0.903</b>	<b>1.122</b>	0.957
<b>Test Set</b>			
Deep3D [10]	1.11	1.41	1.21
DECA [15]	1.09	1.38	1.18
Ours	0.97	1.21	1.02
MICA [61]	<b>0.90</b>	<b>1.11</b>	<b>0.92</b>

### 5.4. REALY 3D Benchmark

We also evaluated our method on the most recent REALY benchmark [14] for single-image 3D face reconstruction from frontal and side view images. Our results, as presented in Tables 2, demonstrate a significant improvement over DECA by **15%** and MICA by **18%**. Our method also stands in the **top 3** (out of 18 methods) on the REALY benchmark challenge outperforming most of the existing methods. The only two methods better than ours are HRN [12] and Deep3D [13]. However, they both generate only frontal mesh, while our method generates a complete head model and is trained in an **unsupervised** setting.

#### 5.4.1 Diversity metrics

One of the important features of a good 3D face model is how diverse it's synthesized meshes are. Similar to prior works that generate 3D meshes [1, 50], we measure global diversity as the mean vertex distance over all possible pairs of  $n$  meshes. Table 3 reports the diversity values for  $n = 1000$  meshes synthesized by MICA[61], DECA[15] and Deep3D[10].

As illustrated in Fig. 6, our model captures head shapes better than DECA model which produces similar-looking head shapes. We observe a significant improvement in diversity statistics over MICA that only predicts a smooth and neutral mesh and was trained on limited 3D data.

Table 2: Single image reconstruction error on REALY Benchmark for Frontal and Side-view images (lower is better).

Method	@nose	@mouth	@forehead	@cheek	@all
<b>Front View</b>					
HRN	1.722	1.357	1.995	1.072	1.537
Deep3D	1.719	1.368	2.015	1.528	1.657
<b>Ours</b>	<b>1.656</b>	<b>2.087</b>	<b>2.102</b>	<b>1.141</b>	<b>1.746</b>
GANFit	1.928	1.812	2.402	1.329	1.868
DECA	1.697	2.516	2.394	1.479	2.010
PRNet	1.923	1.838	2.429	1.863	2.013
EMOCA	1.868	2.679	2.426	1.438	2.103
MICA	1.585	3.478	2.374	1.099	2.134
RingNet	1.934	2.074	2.995	2.028	2.258
<b>Side View</b>					
HRN	1.642	1.285	1.906	1.038	1.468
Deep3D	1.749	1.411	2.074	1.528	1.691
<b>Our</b>	<b>1.576</b>	<b>2.218</b>	<b>2.142</b>	<b>1.112</b>	<b>1.762</b>
PRNet	1.868	1.856	2.445	1.960	2.032
DECA	1.903	2.472	2.423	1.630	2.107
EMOCA	1.867	2.636	2.448	1.548	2.125
MICA	1.525	2.636	2.448	1.548	2.125
RingNet	1.921	1.994	3.081	2.027	2.256

Table 3: Diversity values of randomly generated meshes for various methods. Higher is better.

-	Deep3D [10]	DECA [15]	MICA [61]	Ours
DIV	0.18	1.13	0.21	<b>1.67</b>

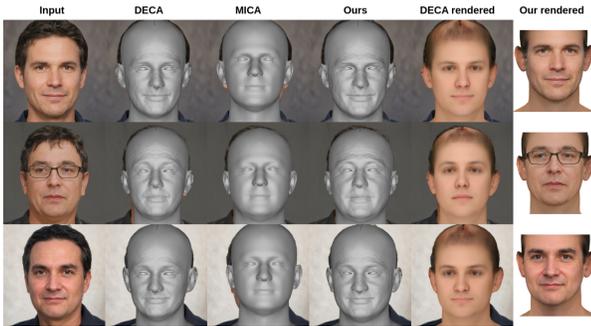


Figure 6: 3D Face Reconstruction comparison with DECA[15] & MICA [61]. DECA is unable to model head shape accurately and provides a smooth texture. MICA only outputs a smooth, neutral and frontal mesh. Our method produces a more detailed mesh with photo-realistic texture, capturing wrinkles, expressions and head shape accurately.

#### 5.4.2 Evaluation of Rendered Faces

We evaluate the quality of texture by comparing qualitatively to a recently proposed method, OSTeC [18]. We also quantitatively compare texture and rendered faces by rendering it on a mesh with other methods, including LiftedGAN [49], DECA [15], and OSTeC [18].

Fig. 7 shows the visual comparison between OSTeC [18] and our model. OSTeC uses latent optimization based GAN inversion which adds random artifacts in the generated image like the beard appearing on some patches of the face. Furthermore, this leads to inconsistency while stitching tex-

tures across different poses.

For quantitative evaluation of our texture and rendered faces, we randomly sampled 1K images from a pretrained StyleGAN2 generator. We use the same set of  $\mathbf{z}$  to generate rendered 3D faces using our method. The corresponding 2D face images are used to perform 3D reconstruction using other methods[15, 18, 49]. To compare how well the texture preserves identity, we measure identity similarity between the input image and the corresponding 3D faces rendered at multiple poses, including the original pose,  $0^\circ$ ,  $\pm 15^\circ$ ,  $\pm 30^\circ$ , and  $\pm 45^\circ$ . The observations can be found in Table ?? . As inferred from the table, our method not only performs better in capturing the identity for the front poses but also for a large number of side poses. More evaluations can be found in supplementary section A.

Method	same pose	$0^\circ$	$\pm 15^\circ$	$\pm 30^\circ$	$\pm 45^\circ$
LiftedGAN [49]	0.90	0.88	0.854	0.83	0.804
DECA [15]	0.869	0.799	0.76	0.69	0.63
OSTeC [18]	0.952	0.939	0.921	0.906	0.88
<b>Ours</b>	<b>0.999</b>	<b>0.995</b>	<b>0.988</b>	<b>0.972</b>	<b>0.941</b>

Table 4: ID similarity comparison between the input image and the corresponding 3D face rendered at various poses.

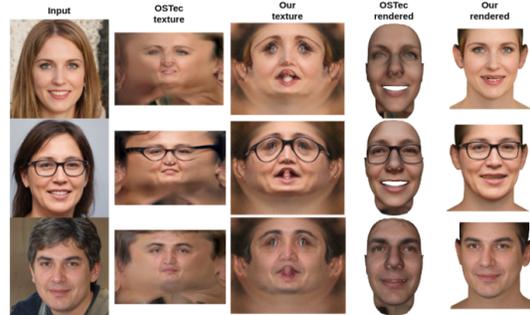


Figure 7: Qualitative comparison of synthesized texture and rendered results with OSTeC [18]. OSTeC produces a smooth texture by stitching multiple images, often leading to artifacts. Our approach can synthesize a better textured mesh outperforming OSTeC in preserving the details.

#### 5.5. 3D face manipulation

Once our end-to-end pipeline is trained and is able to generate albedo and mesh corresponding to a  $\mathbf{z} \in \mathbb{R}^{512}$  or  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ , it opens up an avenue for latent space manipulations in 3D rendered faces. This section shows examples of editing 3D faces directly from the latent space or text.

**Latent space manipulation:** Given a latent code  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ , we can get the modified latent code  $\mathbf{w}_{edit} \in \mathbb{R}^{18 \times 512}$  by  $\mathbf{w}_{edit} = \mathbf{w} + \alpha \mathbf{n}$ , where  $\mathbf{n} \in \mathbb{R}^{512}$  is a vector orthogonal to the semantic boundary and  $\alpha \in \mathbb{R}$  is a constant. More details on latent space manipulations using semantic boundaries can be found in [48, 42]. Fig. 8 shows

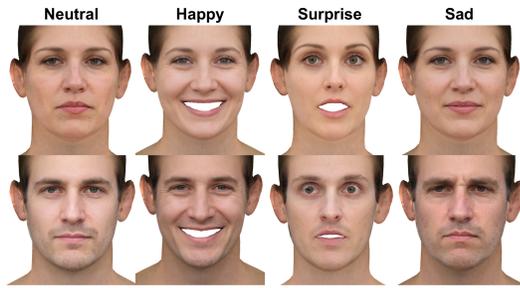


Figure 8: Generating different expressions by manipulating  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$ . The first column in each row is the neutral expression corresponding to randomly sampled  $\mathbf{w}$ . The subsequent columns show different expressions for the same identity by varying  $\mathbf{w}$ .

the result of latent space manipulations for generating different expressions of the same identity.

**Text based 3D Face Editing:** Similar to latent space manipulations, we can perform text-based 3D face editing using Contrastive Language-Image Pretraining (CLIP) models [40]. We used StyleCLIP [38] model to enable text-based editing. The text-guided latent optimization and latent residual mapper strategies of StyleCLIP are implemented in  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$  space and make it easy for us to plug into our pipeline. This paper shows results for the text-guided latent optimization approach, which can be trained for random text queries. Given a  $\mathbf{w}$  for a face, StyleCLIP can produce an updated latent code  $\mathbf{w}' \in \mathbb{R}^{18 \times 512}$  corresponding to a given text query and input  $\mathbf{w}$ . Fig. 1(c) illustrates some sample results generated for the text-based 3D face editing. Refer to supplementary Section B.5 for more results on text-based editing.

## 6. Conclusion and Future Work

In this paper, we attempt to develop a high-quality 3D face generation pipeline. We propose AlbedoGAN that synthesizes albedo and generalizes well over multiple poses capturing intrinsic details of the face. Our approach generates meshes that capture high-frequency details like face wrinkles. Comprehensive experiments demonstrate superiority of our method over others in predicting detailed mesh and preserving the identity in reconstructed 3D faces. As a consequence of using StyleGAN2 based pipeline, we bring style editing, semantic face manipulations, and text-based editing in 3D faces.

While our pipeline can generate high-quality 3D faces from StyleGAN2's latent space, some issues still need to be addressed. Mesh-based representations are unable to model details like hair. We foresee exploiting topologically uniform mesh, and a NeRF-based approach should be able to capture such facial features. We'll extend our work to incorporate more complex illumination models.

## References

- [1] Victoria Fernández Abrevaya, Adnane Boukhayma, Stefanie Wuhler, and Edmond Boyer. A decoupled 3d facial shape model by adversarial training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9419–9428, 2019. 2, 7
- [2] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Only a matter of style: Age transformation using a style-based regression model. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021. 5
- [3] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces, 1999. 2, 3
- [4] Chen Cao, Yanlin Weng, Shun Zhou, Yiyang Tong, and Kun Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2013. 3, 7
- [5] Eric R. Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. Pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5799–5809, June 2021. 2
- [6] Hang Dai, Nick Pears, William Smith, and Christian Duncan. Statistical modeling of craniofacial shape and texture. *International Journal of Computer Vision*, 128(2):547–571, 2020. 3, 7
- [7] Jiankang Deng, Shiyang Cheng, Niannan Xue, Yuxiang Zhou, and Stefanos Zafeiriou. Uv-gan: Adversarial facial uv map completion for pose-invariant face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7093–7102, 2018. 2
- [8] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 3
- [9] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. 5
- [10] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *IEEE Computer Vision and Pattern Recognition Workshops*, 2019. 2, 3, 4, 7, 8
- [11] Epic Games. Unreal engine. 2
- [12] Biwen et al. A hierarchical representation network for accurate and detailed face reconstruction from in-the-wild images. In *CVPR*, 2023. 7
- [13] Yu et al. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *CVPRW*, 2019. 7
- [14] Zenghao et al. Realy: Rethinking the evaluation of 3d face reconstruction. In *ECCV*, 2022. 7
- [15] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3D face model from in-the-wild images. In *ACM Transactions on Graphics, (Proc. SIGGRAPH)*, volume 40, 2021. 2, 3, 4, 7, 8

- [16] Zhen-Hua Feng, Patrik Huber, Josef Kittler, Peter Hancock, Xiao-Jun Wu, Qijun Zhao, Paul Koppen, and Matthias Rätzsch. Evaluation of dense 3d reconstruction from 2d face images in the wild. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 780–786. IEEE, 2018. 3, 7
- [17] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 2
- [18] Baris Gecer, Jiankang Deng, and Stefanos Zafeiriou. Ostedec: One-shot texture completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7628–7638, June 2021. 2, 3, 8
- [19] Baris Gecer, Alexandros Lattas, Stylianos Ploumpis, Jiankang Deng, Athanasios Papaioannou, Stylianos Moschoglou, and Stefanos Zafeiriou. Synthesizing coupled 3d face modalities by trunk-branch generative adversarial networks. In *European conference on computer vision*, pages 415–433. Springer, 2020. 2
- [20] Baris Gecer, Alexander Lattas, Stylianos Ploumpis, Jiankang Deng, Athanasios Papaioannou, Stylianos Moschoglou, and Stefanos Zafeiriou. Synthesizing coupled 3d face modalities by trunk-branch generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV)*, 2020. 2, 3
- [21] Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [22] Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1155–1164, 2019. 2
- [23] Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. Fast-ganfit: Generative adversarial network for high fidelity 3d face reconstruction. *arXiv preprint arXiv:2105.07474*, 2021. 2
- [24] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18653–18664, 2022. 2
- [25] Jia Guo, Jiankang Deng, Alexandros Lattas, and Stefanos Zafeiriou. Sample and computation redistribution for efficient face detection. *arXiv preprint arXiv:2105.04714*, 2021. 6
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 5
- [27] Alexandru Eugen Ichim, Sofien Bouaziz, and Mark Pauly. Dynamic 3d avatar creation from hand-held video input. *ACM Transactions on Graphics (ToG)*, 34(4):1–14, 2015. 2
- [28] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 2
- [29] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 6
- [30] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2
- [31] Jongyoo Kim, Jiaolong Yang, and Xin Tong. Learning high-fidelity face texture completion without complete face texture. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13990–13999, 2021. 2, 3
- [32] Chunlu Li, Andreas Morel-Forster, Thomas Vetter, Bernhard Egger, and Adam Kortylewski. To fit or not to fit: Model-based face reconstruction and occlusion segmentation from weak supervision. In *arXiv preprint arXiv:2106.09614*, 2021. 2, 3
- [33] 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)*, 36(6), 2017. 2, 3, 4, 5
- [34] Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. Deep appearance models for face rendering. *ACM Transactions on Graphics (ToG)*, 37(4):1–13, 2018. 2
- [35] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *ICML*, 2022. 1
- [36] Nicolas Olivier, Kelian Baert, Fabien Danieau, Franck Multon, and Quentin Avril. Facetunegan: Face autoencoder for convolutional expression transfer using neural generative adversarial networks. *arXiv preprint arXiv:2112.00532*, 2021. 2
- [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 6
- [38] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 9
- [39] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance*, pages 296–301, 2009. 2, 3

- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [2](#), [9](#)
- [41] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. [1](#)
- [42] Aashish Rai, Clara Ducher, and Jeremy R Cooperstock. Improved attribute manipulation in the latent space of stylegan for semantic face editing. In *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 38–43. IEEE, 2021. [6](#), [8](#)
- [43] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. [1](#)
- [44] Daniel Rebaïn, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. Lolnerf: Learn from one look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1558–1567, 2022. [2](#)
- [45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 2022. [1](#)
- [46] Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael Black. Learning to regress 3d face shape and expression from an image without 3d supervision. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#), [3](#), [7](#)
- [47] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. [5](#)
- [48] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. In *TPAMI*, 2020. [6](#), [8](#)
- [49] Yichun Shi, Divyansh Aggarwal, and Anil K Jain. Lifting 2d stylegan for 3d-aware face generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6258–6266, 2021. [8](#)
- [50] Fariborz Taherkhani, Aashish Rai, Quankai Gao, Shaunak Srivastava, Xuanbai Chen, Fernando de la Torre, Steven Song, Aayush Prakash, and Daeil Kim. Controllable 3d generative adversarial face model via disentangling shape and appearance. *arXiv preprint arXiv:2208.14263*, 2022. [2](#), [7](#)
- [51] Ayush Tewari, Michael Zollhöfer, Hyeonwoo Kim, Pablo Garrido, Florian Bernard, Patrick Pérez, and Christian Theobalt. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. *CoRR*, abs/1703.10580, 2017. [2](#), [3](#)
- [52] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. [2](#)
- [53] Mukhiddin Toshpulatov, Wookey Lee, and Suan Lee. Generative adversarial networks and their application to 3d face generation: a survey. *Image and Vision Computing*, 108:104119, 2021. [2](#)
- [54] Unity Technologies. Unity. [2](#)
- [55] Thibaut Weise, Sofien Bouaziz, Hao Li, and Mark Pauly. Realtime performance-based facial animation. *ACM transactions on graphics (TOG)*, 30(4):1–10, 2011. [2](#)
- [56] Xiangnan Yin, Di Huang, Zehua Fu, Yunhong Wang, and Liming Chen. Weakly-supervised photo-realistic texture generation for 3d face reconstruction. *arXiv preprint arXiv:2106.08148*, 2021. [2](#)
- [57] Xi Yin, Ying Tai, Yuge Huang, and Xiaoming Liu. Fan: Feature adaptation network for surveillance face recognition and normalization. In *Proceedings of the Asian Conference on Computer Vision*, 2020. [5](#), [6](#)
- [58] Yin Yu, Ghasedi Kamran, Wu HsiangTao, Yang Jiaolong, Tong Xi, and Fu Yun. Expanding the latent space of stylegan for real face editing. *arXiv preprint arXiv:2204.12530*, 2022. [5](#)
- [59] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C Bühler, Xu Chen, Michael J Black, and Otmar Hilliges. Im avatar: Implicit morphable head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13545–13555, 2022. [2](#)
- [60] Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. Mofan-erf: Morphable facial neural radiance field. *arXiv preprint arXiv:2112.02308*, 2021. [2](#)
- [61] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Towards metrical reconstruction of human faces. In *European Conference on Computer Vision*, 2022. [2](#), [3](#), [6](#), [7](#), [8](#)