

Learning Person-Specific Animatable Face Models from In-the-Wild Images via a Shared Base Model

Anonymous CVPR submission

Paper ID 10007

Abstract

001 *Training a generic 3D face reconstruction model in a self-*
 002 *supervised manner using large-scale, in-the-wild 2D face*
 003 *image datasets enhances robustness to varying lighting con-*
 004 *ditions and occlusions while allowing the model to cap-*
 005 *ture animatable wrinkle details across diverse facial expres-*
 006 *sions. However, a generic model often fails to adequately*
 007 *represent the unique characteristics of specific individuals.*
 008 *In this paper, we propose a method to train a generic base*
 009 *model and then transfer it to yield person-specific models by*
 010 *integrating lightweight adapters within the large-parameter*
 011 *ViT-MAE base model. These person-specific models excel*
 012 *at capturing individual facial shapes and detailed fea-*
 013 *tures while preserving the robustness and prior knowledge*
 014 *of detail variations from the base model. During training,*
 015 *we introduce a silhouette vertex re-projection loss to*
 016 *address boundary “landmark marching” issues on the 3D*
 017 *face caused by pose variations. Additionally, we employ*
 018 *an innovative teacher-student loss to leverage the inher-*
 019 *ent strengths of UNet in feature boundary localization for*
 020 *training our detail MAE. Quantitative and qualitative ex-*
 021 *periments demonstrate that our approach achieves state-of-*
 022 *the-art performance in face alignment, detail accuracy, and*
 023 *richness. The code will be released to the public upon the*
 024 *acceptance of this paper.*

025 1. Introduction

026 The reconstruction of 3D faces from 2D images has gar-
 027 nered considerable attention recently [5, 8, 12, 18, 20, 42,
 028 60, 62, 72, 76], with applications spanning diverse fields
 029 such as 3D avatar creation [1, 24, 30], face recognition [2,
 030 3, 50], and face animation driven by speech [11, 17, 47, 48]
 031 or video [27, 47, 81]. Leveraging deep learning, the major-
 032 ity of recent methods [5, 6, 8, 12, 15, 18, 20, 21, 26, 28,
 033 34, 38, 42, 49, 56–59, 62, 66–70, 72, 76, 79, 80] focus on
 034 reconstructing 3D faces from in-the-wild images, primarily
 035 employ a unified set of model weights across images of dif-

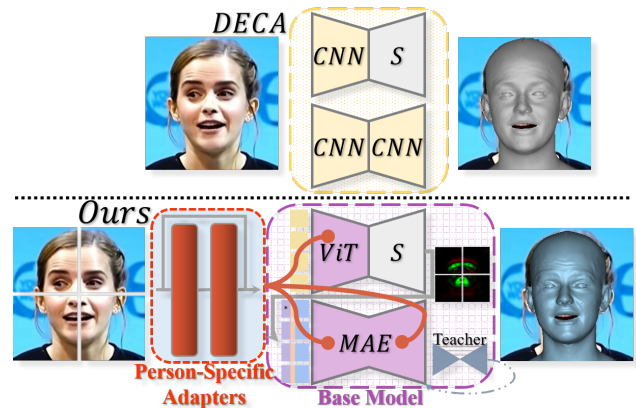


Figure 1. In contrast to previous work (e.g., DECA), we first develop a scalable, high-capacity base model (purple) and then transfer it to create person-specific models by integrating lightweight, person-specific adapters (red).

036 ferent individuals. While robust, these methods often under-
 037 fit individual-specific features. However, in real-world
 038 scenarios, multiple images or a video of the same person are
 039 often available, allowing models to focus on reconstructing
 040 that specific individual more accurately—a context where
 041 current methods still have limitations.

042 In contrast to previous work, we first develop a scal-
 043 able, high-capacity base model, trained in a self-supervised
 044 manner on extensive 2D images, and then transfer it to
 045 yield person-specific models. Our base model reconstructs
 046 feature-aligned 3D faces from in-the-wild images in real
 047 time, capturing expression-related wrinkle details (animat-
 048 able features). When multiple images or a video of an in-
 049 dividual are available, we can transfer the generic model
 050 by integrating a small set of person-specific parameters,
 051 achieving more precise real-time reconstruction (30 fps on
 052 an Nvidia GeForce RTX 3090) while maintaining robust-
 053 ness against occlusions and preserving priors of various fa-
 054 cial details associated with expressions.

055 Vision Transformer (ViT) [16] and Masked Autoencoder
 056 (MAE) [22] are models with strong expressive capabilities,

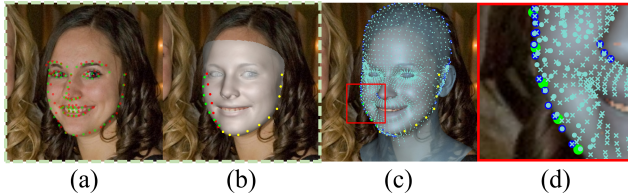


Figure 2. **Motivation of our silhouette vertex re-projection loss.** (a) 2D landmark ground truth [55] (green) and projected 3D landmarks provided by 3DDFA-v2 [20] (red). (b) The occluded boundary 3D landmarks (red) “move” to the face silhouette in the 2D annotation (green). (c) All vertices of the coarse model reconstructed by our method (light cyan, with crosses indicating normals pointing inward towards the image and dots indicating normals pointing outward), and the outer edge points of the 3D face (blue). (d) A zoom-in of the red box area from (c).

effectively gathering both local and global dependencies in visual data. They also excel in achieving robust performance and generalization on target tasks with limited training data by leveraging pre-training on large-scale datasets, aligning well with our design. Capitalizing on these advantages, We propose a ViT-MAE architecture for generic 3D face reconstruction, leveraging differentiable rendering [46] for self-supervised training. The model learns a parametric face model for coarse reconstruction, refined by a UV displacement map. This trained generic model serves as a base model, which can be further transferred to a person-specific model by integrating lightweight adapter modules [23] within the transformer layer. Our approach enables the single-image face reconstruction model to process multiple images or videos, fully leveraging the data to refine the reconstruction outcomes.

When a 3D face is projected onto an image, inner and boundary landmarks outlining facial features and the cheek are also projected. Face reconstruction methods [26, 56, 57, 66–68, 79] constrain the reconstructed 3D face by minimizing the error between the projected 3D landmarks and the annotated ground truth. However, in non-frontal views, some landmarks, especially occluded boundary ones, become invisible (red points in Fig. 2.(b)), making accurate annotation difficult. In 2D landmark annotation, boundary landmarks for the 3D cheek “shift” to align with the face silhouette, causing misalignment—known as “landmark marching” [82]. Our innovative silhouette vertex re-projection loss addresses this by aligning 2D silhouette landmarks with silhouette edge vertices on the 3D model based on current vertex normal distribution (Fig. 2.(c)&(d)). Additionally, using dense silhouette edge vertices as loss candidates enhances the model’s sensitivity to normal errors, given the high variance in manual silhouette landmark annotations [55] and their tendency to spread along the boundary tangent [25]. Experimental results show that our method surpasses prior approaches in face alignment.

We define facial details as functions of identity and “ten-

sion” [44] within facial geometry, which vary with expressions. The local correspondence in UV space between the detailed displacement map and the unwrapped image texture makes a UNet [51] well-suited for learning it, with skip connections aiding in the precise localization of feature boundaries. However, the skip connections indiscriminately convey all facial details to the displacement map, hindering the modeling of the detail changes due to facial deformation, and also making it challenging to fill in the “holes” caused by occlusion. We employ a masked autoencoder (MAE) [22] for detail generation, using a consistency loss [18] to decouple identity and “tension”-related details. During training, a UNet acts as a “teacher” to guide the MAE to learn detailed feature boundaries accurately. Integrating the strengths of both UNet and MAE in detail recovery and animation, our model outperforms previous works in wrinkle diversity and accuracy without relying on 3D data. Our ablation study underscores the effectiveness.

In summary, our main contributions are:

Personalized Face Model. To the best of our knowledge, we are the first to construct a person-specific 3D face reconstruction model by transferring a large-scale generic model. The person-specific adapters within the generic ViT-MAE enhance the reconstruction from that person’s images on both coarse and fine scales.

Silhouette Vertex Re-projection Loss. We introduce an innovative loss function that aligns the annotated 2D silhouette landmarks with the 3D facial boundary edges, providing more robust facial contour constraints.

Animatable Details and Teacher-Student Architecture. We innovatively employ a teacher-student architecture that leverages UNet’s inherent strength in feature boundary localization alongside MAE’s animation and robustness capabilities, enabling the MAE, which is ultimately used for reconstruction and animation, to capture both identity-related and tension-related details effectively.

2. Related Work

Monocular Coarse Reconstruction. Landmark loss is widely used in self-supervised approaches [26, 37, 56, 57, 66–68, 79] to ensure the alignment of projected 3D landmarks with 2D landmarks detected by face alignment techniques. However, the “landmark marching” phenomenon remains a persistent challenge. [56, 68] use dynamic boundary landmarks determined by head pose to maintain correct correspondence but overlook the influences of facial shape and expression. [26] defines silhouette landmarks on horizontal mesh lines, requiring manual redefinition for new topologies. [37] renders a projection area to identify occluded boundary landmarks, though it is time-intensive.

Animatable Detail Reconstruction. Adding facial details greatly improves model authenticity and expressiveness. While many studies have achieved high-quality detail

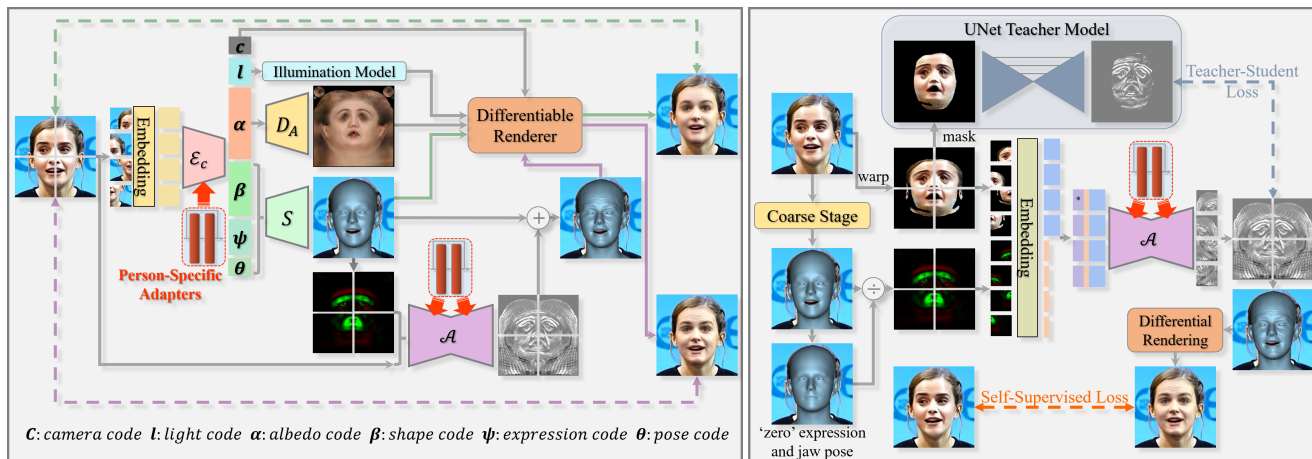


Figure 3. **Illustration of our overall architecture.** **Left box:** End-to-end learning framework of our coarse and detail stage. Given an image, we first regress the FLAME parameters to obtain a coarse facial shape (following the green arrows) with the help of our proposed silhouette vertex re-projection loss. We then use an MAE to reconstruct animatable facial details from the warped image texture and current facial tension [44] (following the purple arrows). During person-specific transfer, we integrate lightweight, trainable adapters (red) while freezing the base model weights. **Right box:** The training pipeline for animatable details. We first train a UNet teacher model that can effectively recover visible details from the current perspective. Under the guidance of the UNet, we train our animatable detail MAE \mathcal{A} using a teacher-student paradigm.

147 generation [8, 21, 26, 34, 38, 49, 58, 59, 69, 72, 80], animat- 177
 148 able details remain underexplored. Such details are essen- 178
 149 tial for lifelike avatars that respond naturally to expres-
 150 sions. Methods like [31, 78] generate impressive animatable
 151 details from textured meshes or neutral faces but struggle to
 152 differentiate static from dynamic features. While [14] and
 153 [32] generate animatable details, neither method supports
 154 3D face reconstruction from 2D images. [36] uses Style-
 155 GAN2 to animate details from facial images or 3D meshes
 156 but lacks robustness in varied lighting. Closest to our ap-
 157 proach, [5] and [18] reconstruct 3D faces with animatable
 158 details from single images; however, [5] is limited by its de-
 159 pendency on synthetic data for supervised training, and [18]
 160 achieves reasonable results but faces challenges in detail ac-
 161 curacy and the realism of expression-related variations.

162 **Multiple Images or Monocular Video.** When multiple
 163 images or videos of a subject are available, we aim to fully
 164 utilize this data. Optimization-based methods [41, 53, 54,
 165 65] often face limitations such as slow inference, difficulty
 166 adapting to unseen facial areas, and insufficient geometric
 167 detail. Learning-based methods [15, 21, 38, 45, 60, 77]
 168 leverage deep networks to integrate different viewpoints but
 169 struggle to capture intricate facial details. While [21] and
 170 [38] can reconstruct high-detail textured geometry from sin-
 171 gle videos, they fail to model unique expression-specific
 172 features or wrinkles absent in the video. Critically, none
 173 of these methods produce animatable facial details.

174 3. Method

175 Our base model and transferred person-specific models
 176 share a two-stage self-supervised training framework. This

section details the coarse and detailed stages, followed by 177
 the transfer process to obtain a person-specific model. 178

179 3.1. Preliminary

180 **3D Geometry Model.** FLAME [33] is a statistical 3D head 180
 181 model that, given identity shape parameters $\beta \in \mathbb{R}^{|\beta|}$, fa- 181
 182 cial expression parameters $\psi \in \mathbb{R}^{|\psi|}$, and pose parameters 182
 183 $\theta \in \mathbb{R}^{3k+3}$ for rotations around $k = 4$ joints (neck, jaw, and 183
 184 eyeballs) and global rotation, outputs a mesh $S(\beta, \theta, \psi)$ 184
 185 with $n_v = 5023$ vertices and $n_f = 9976$ faces. 185

186 **Appearance Model.** We use a texture statistical model 186
 187 that aligns the Basel Face Model’s albedo space [40] with 187
 188 FLAME’s UV layout [18], which outputs a FLAME tex- 188
 189 ture map $A(\alpha) \in \mathbb{R}^{d \times d \times 3}$, where $d = 256$, given texture 189
 190 parameters $\alpha \in \mathbb{R}^{|\alpha|}$. 190

191 **Camera Model.** The camera model is parameterized by 191
 192 $\mathbf{c} = (s, \mathbf{t})$. An orthographic projection transformation is 192
 193 used to project 3D mesh vertices into the image space, for- 193
 194 mulated as $\mathbf{u} = s\Pi(\mathbf{v}) + \mathbf{t}$, where \mathbf{v} is a vertex in the 194
 195 3D mesh, $\Pi \in \mathbb{R}^{2 \times 3}$ is the orthographic projection matrix, 195
 196 $s \in \mathbb{R}$ is the isotropic scale, and $\mathbf{t} \in \mathbb{R}^2$ is the 2D transla- 196
 197 tion. 197

198 **Illumination Model.** The spherical harmonics illu- 198
 199 mination model [43] is adopted to estimate the illumina- 199
 200 tion conditions in the input image. The shaded face 200
 201 texture is computed as $U = U(A, \mathbf{l}, N)$ with $U_{i,j} =$ 201
 202 $A_{i,j} \odot \mathcal{H}_{i,j} = A_{i,j} \odot \sum_{k=1}^9 \mathbf{l}_k H_k(N_{i,j})$, where \odot denotes 202
 203 the Hadamard product, $A_{i,j}$, $N_{i,j}$, and $U_{i,j}$ represent the 203
 204 albedo, surface normal, and shaded texture in UV coordi- 204
 205 nates respectively, $H : \mathbb{R}^3 \rightarrow \mathbb{R}$ are the Spherical Harmon- 205
 206 ics (SH) basis functions, and $\mathbf{l} \in \mathbb{R}^9$ is the SH coefficient. 206

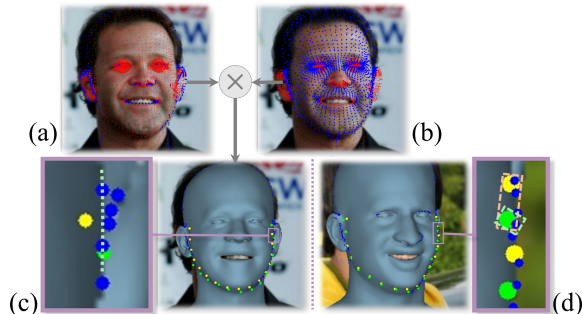


Figure 4. **Our silhouette vertex re-projection loss construction method and its superiority.** (a) Edge points (blue) selected from all model vertices (red) based on vertex normals. (b) Candidate region (blue) for boundary edge vertices excludes the nose and the ear regions. (c) & (d) Two typical scenarios.

3.2. Generic Coarse Reconstruction

We start by employing self-supervised learning to develop a base model that robustly reconstructs any in-the-wild image with fine details. Our base model first learns a coarse reconstruction. A vision transformer [16] \mathcal{E}_c serves as the coarse encoder. It splits a 2D face image I into 16×16 patches, embeds each patch linearly and adds positional embeddings to form a sequence of tokens. A learnable “classification token” ($\mathbf{z}_0^0 = \mathbf{x}_{\text{class}}$) is prepended to this sequence, which is subsequently passed through L transformer blocks. The output classification token \mathbf{z}_L^0 is passed through an MLP with one hidden layer to generate a latent code comprising FLAME parameters β , θ , and ψ , along with albedo parameters α , camera parameters \mathbf{c} , and lighting parameters \mathbf{l} . With FLAME and albedo parameters, we generate a textured 3D mesh, which, combined with camera and lighting parameters, enables differentiable rendering to produce the reconstructed facial image I_r (Fig. 3, left, green arrows).

Silhouette Vertex Re-Projection Loss. The first 17 landmarks $\mathbf{k}_i \in \mathbb{R}^2, i \in 1, \dots, 17$ from the annotated ground truth are 2D silhouette landmarks that should align with the outer edge of the 3D face from the current view. Our silhouette vertex re-projection loss naturally enforces this alignment. Compared to FLAME’s dynamic landmark marching approach [33], which may misalign landmarks under variations in facial shape and expression, our method, illustrated in Fig. 4.(c)&(d), achieves more accurate matching in two typical scenarios: (1) when the dynamic boundary landmark (yellow) provided by the FLAME model does not accurately locate on the model’s edge (green dashed line); and (2) when the FLAME algorithm necessitates matching the ground truth landmark (green) with the yellow point as indicated by the orange box, whereas our method matches it to the blue edge point shown in the green box, providing a more reasonable alignment.

We define the ‘zero-pose’ boundary landmarks of the FLAME model $\tilde{\mathbf{k}}_i$ as those selected by FLAME when given

a zero-pose input. If a landmark $\tilde{\mathbf{k}}_i$ is non-occluded, it is constrained using the vanilla landmark re-projection; if occluded, we match it to the nearest edge point in V_e and calculate the \mathbb{L}_1 loss. Occlusion is determined by vertex normals: if the z-direction points outward, it is non-occluded; otherwise, it is occluded. The determination of edge point set V_e is as follows: if the normals of two vertices at the ends of an edge have opposite signs in the z-direction, it is considered that both of these vertices are located on the edge of the 3D model, we then project them onto the image plane (Fig. 4.(a)). Furthermore, as shown in Fig. 4.(b), we exclude the edges of the nose and ear regions using pre-segmented vertex labels $\mathbf{Z} \in \{0, 1\}^{5023}$. We refer to the cleaned edge points as the 3D “silhouette vertices”.

The silhouette vertex re-projection loss is defined as follows:

$$\mathcal{L}_{sil} = \sum_{i, N_{z,i} > 0} d(\mathbf{k}_i, \tilde{\mathbf{k}}_i) + \sum_{i, N_{z,i} < 0} \min_{\mathbf{v}_e \in \tilde{V}_e} d(\mathbf{k}_i, \mathbf{v}_e), \quad (1)$$

where $N_{z,i} = (N_{\tilde{\mathbf{k}}_i})_z$, $\tilde{\mathbf{k}}_i, i \in 1, \dots, 17$ represents the 3D landmark coordinates that match the 2D landmark \mathbf{k}_i in the zero-pose, N represents the vertex normal, $(\bullet)_z$ denotes the z-component. $\tilde{V}_e = V_e \odot \tilde{\mathbf{Z}} \in \mathbb{R}^{5023 \times 2}$, $\tilde{\mathbf{Z}} \in \{0, 1\}^{5023 \times 2}$. If $Z_i = 0$, then $\tilde{\mathbf{Z}}_i = (0, 0)^T$; otherwise, $\tilde{\mathbf{Z}}_i = (1, 1)^T$.

Overall Losses for Coarse Reconstruction. The base model learns coarse reconstruction using the total loss \mathcal{L}_C :

$$\mathcal{L}_C = \mathcal{L}_{sil} + \mathcal{L}_{inL} + \mathcal{L}_{spL} + \mathcal{L}_{pho} + \mathcal{L}_{per} + \mathcal{L}_{reg} + \mathcal{L}_{sc}, \quad (2)$$

where \mathcal{L}_{inL} is the landmark loss for static inner landmarks, \mathcal{L}_{spL} is the special landmark pairs loss that is calculated on a set of landmark pairs (e.g., upper/lower eyelid or lip landmarks) to constrain features like eye and mouth opening in a translation-invariant manner. \mathcal{L}_{pho} is the photometric loss commonly used in self-supervised methods [15, 18]. $\mathcal{L}_{per} = \mathcal{L}_{id} + \mathcal{L}_{emo} + \mathcal{L}_{lr}$ combines three perceptual losses to ensure high-level identity [18] and emotion consistency [12], as well as accurate lip movements [19]. \mathcal{L}_{reg} includes regularization losses for β , ψ , and α . Finally, \mathcal{L}_{sc} is the shape consistency loss [18]

3.3. Generic Detail Reconstruction

Our base model subsequently learns a displacement map $D \in [-0.01, 0.01]^{d \times d}$ to refine the FLAME geometry. We use an MAE [22] to capture animatable facial details (Fig. 3, right), defined as $D = \mathcal{A}(I, T_{UV}, \mathbf{x}_a)$, where I is the input image, T_{UV} is the current face tension map, and \mathbf{x}_a is a learnable, patch-wise detail latent that is shared across all individuals. Unlike the vanilla MAE, we do not apply masking (i.e., use a zero masking ratio) because our task differs significantly from standard MAE applications, requiring the autoencoder to leverage all available information to accurately recover facial details aligned with the input images.

Facial Tension. Facial tension quantifies the vertex-wise compression or expansion on the 3D mesh caused by deformation from a neutral expression with a closed mouth to the current expression [5, 44]. We propose a new method for calculating the tension at vertex v_i :

$$t_{v_i}(S, S') = \frac{1}{E_i} \sum_{E_i}^{j=1} \frac{e^{k\|e'_j\|} - e^{k\|e_j\|}}{e^{k\|e'_j\|} + e^{k\|e_j\|}}, \quad (3)$$

where e_1, \dots, e_{E_i} are the E_i edges connected to v_i in $S = S(\beta, \theta, \psi)$, and e'_1, \dots, e'_{E_i} are the E_i edges connected to v'_i in $S' = S(\beta, \mathbf{0}, \mathbf{0})$. Here, $\|\bullet\|$ denotes the edge length, and k is a fixed scaling factor. Compression in the vertex neighborhood results in positive tension, while stretching yields negative tension. Our tension metric $t_{v_i}(S, S')$ satisfies: 1) Antisymmetry: $t_{v_i}(S', S) = -t_{v_i}(S, S')$; and 2) Boundedness: $\forall S, S', -1 < t_{v_i}(S, S') < 1$, making this tension calculation more suitable as input to a neural network. The tension of S can be represented as $\{t_{v_i}(S, S')\}$. Using the mapping relationship of UV coordinates from the FLAME mesh S , we derive the tension map $T_{UV}(\beta, \theta, \psi)$ in UV space.

Teacher-Student Strategy. Unwrapping the input image to UV space using the reconstructed coarse FLAME geometry creates a local correspondence between the image texture and the UV space displacement map. Thus, a UNet [51] is more suitable for estimating the displacement map from the UV image texture than an encoder-decoder or autoencoder, as its skip connections efficiently utilize local input information. However, the skip connections make it difficult to animate facial details according to facial deformations, and non-frontal poses or occlusions may cause incompleteness in the unwrapped texture, leading to missing details in the reconstruction. Autoencoder structures, on the other hand, have unique advantages in terms of animation and robustness to occlusion. Therefore, we train a shallow UNet detail reconstruction network as a teacher model, using its estimates as pseudo ground truth to guide our autoencoder in detail reconstruction. This teacher model offers more direct supervision compared to shape-from-shading photometric loss, as 3D facial details and 2D shading on rendered images do not have a one-to-one mapping—multiple ways exist to add details that result in the same shadows in the rendered output.

UNet Teacher Training. We train the UNet teacher network \mathcal{D}_{UNet} for detail estimation by minimizing:

$$\mathcal{L}_{UNet} = \mathcal{L}_{phoD} + \mathcal{L}_{mrf} + \mathcal{L}_{smo} + \mathcal{L}_{regD}. \quad (4)$$

where \mathcal{L}_{pho} is the photometric loss for detail rendering [18]. \mathcal{L}_{mrf} is an ID-MRF loss [74], computed on the *conv3.2* and *conv4.2* layers of VGG19 [63], encouraging the model to capture high-frequency geometric details. \mathcal{L}_{smo} is a

smoothness loss prevents overly sharp or high-frequency artifacts in the reconstructed details. \mathcal{L}_{regD} is the detail regularization loss regularizes the estimated displacements to reduce noise and artifacts. For more details, please refer to the supplementary materials.

Teacher-Student Loss. With the UNet teacher network trained, we can employ the teacher-student loss to aid the training of our detail MAE \mathcal{A} . The teacher-student loss measures the similarity between the displacement given by the pretrained teacher network UNet and the student network \mathcal{A} . The photometric loss and SSIM [75] loss are used:

$$\mathcal{L}_{Tchr} = \mathcal{L}_{phoD}(D, D_{UNet}) + \mathcal{L}_{ssim}(D, D_{UNet}). \quad (5)$$

Overall Losses for Animatable Details. We train \mathcal{A} on video datasets [35, 71], where each mini-batch contains images from different frames of the same video. This allows the model to learn animatable facial details that vary with changes in facial tension. In total, we optimize:

$$\mathcal{L}_{animD} = \mathcal{L}_{UNet} + \mathcal{L}_{Tchr} + \mathcal{L}_{sym} + \mathcal{L}_{dc}, \quad (6)$$

where $\mathcal{L}_{UNet} = \mathcal{L}_{phoD} + \mathcal{L}_{mrf} + \mathcal{L}_{smo} + \mathcal{L}_{regD}$ includes the losses used for training the UNet teacher model, which are also employed during the training of the student model \mathcal{A} . \mathcal{L}_{sym} is the soft symmetry loss employed to enhance the model’s robustness in occlusion regions and reduce boundary artifacts. \mathcal{L}_{dc} is the detail consistency loss:

$$\mathcal{L}_{dc} = \mathcal{L}_{animD}(I, \mathcal{D}_{UNet}(I), \mathcal{D}(I', T_{UV})), \quad (7)$$

where, ensures that for images of the same individual, swapping the input images makes no difference, as they should convey the same identity information. Details from I' and T_{UV} should be consistent with image I and the pseudo ground truth $\mathcal{D}_{UNet}(I)$ from the teacher UNet.

3.4. Person-Specific Transfer

Given multiple images or videos of a person, we can transfer the base model to yield a person-specific model. The geometry of the person-specific model better aligns with the shapes and boundaries of the inputs, capturing the unique facial details of the individual. The person-specific model retains the base model’s robustness to pose and occlusion, as well as valuable priors regarding dynamic details.

We achieve this transfer by incorporating lightweight modules, δ_{PS} , known as “adapters” [23], between layers of the base model (Fig. 3). The base model parameters remain fixed, and only the adapter parameters are trained. Due to the residual structure of the adapters, the modifications to the base model are incremental. Positioned within the transformer blocks, adapters δ_{PS} are added after the feed-forward layer, preceding layer normalization. Additionally, the learnable patch-wise detail latent \mathbf{x}_d and the layer normalization parameters are also trained. More details are provided in the supplementary materials.

Table 1. Reconstruction error across different datasets.

Method	300-W [55]			300-VW [61]			FaceScape [78]	
	boundary↓	inner↓	overall↓	boundary↓	inner↓	overall↓	inner error↓	overall error↓
DECA [18]	0.0576	0.0402	0.0446	0.0563	0.0426	0.0460	0.0477 (± 0.0097)	0.0496 (± 0.0103)
EMOCA [12]	0.0579	0.0467	0.0495	-	-	-	-	-
EMOCA-v2 [12, 19]	0.0590	0.0377	0.0430	0.0553	0.0507	0.0519	0.0585 (± 0.0111)	0.0600 (± 0.0109)
SynergyNet [76]	-	0.0545	0.0545	-	0.0651	0.0651	0.0589 (± 0.0117)	0.0785 (± 0.0230)
3DDFA-v2 [20]	-	0.0470	0.0470	-	0.0524	0.0524	0.0514 (± 0.0104)	0.0700 (± 0.0218)
Ours-base (w/o \mathcal{L}_{sil}) [†]	0.0616	0.0401	0.0455	0.0660	0.0554	0.0580	-	-
Ours-base	0.0559	0.0348	0.0401	0.0585	0.0456	0.0488	0.0404 (± 0.0126)	0.0429 (± 0.0128)
Ours-base w/ δ_{PS}	-	-	-	0.0359	0.0217	0.0252	0.0253 (± 0.0045)	0.0305 (± 0.0060)

* A new version of EMOCA that incorporates perceptual lip reading loss [19] and produces better lip and eye alignment compared to the original model [12].

[†] Our base model trained with DECA’s landmark re-projection loss [18], without the proposed silhouette vertex re-projection loss (Eqn. 1).

389 Person-specific transfer is also divided into a coarse
390 stage and a detailed stage. Given multiple images or video
391 frames of an individual, the model is trained using the loss
392 function \mathcal{L}_C (Eqn. 2) in the coarse stage and teacher su-
393 pervision loss \mathcal{L}_{Tchr} (Eqn. 5) in the detail stage. On an
394 NVIDIA GeForce RTX 3090, several minutes of training
395 yields significantly improved reconstruction results. Spe-
396 cially, since the occluded areas are not optimized by the
397 losses during transfer, they retain the priors from the base
398 model, allowing for reasonable generation of detailed fa-
399 cial features based on the extensive face details learned by
400 the base model. Moreover, due to the strong identity con-
401 sistency of the coarse shape of the person-specific model
402 across different frames, variations in shape between frames
403 are more attributed to expression changes. This facilitates
404 better decoupling of identity and expression parameters in
405 the facial statistical model, leading to more precise ex-
406 pression estimation and dynamic, expression-related details
407 specific to the individual.

408 4. Experiments

409 4.1. Implementation Details

410 **Datasets.** We train our model using the publicly available
411 datasets BUPT-Balancedface [73], Celeb-DF (v2) [35], and
412 MEAD [71]. For each image, facial landmarks are automat-
413 ically annotated using HRNet [64], and unreliable images
414 are filtered based on the estimated per-landmark heatmaps.
415 We utilize a facial skin region segmentation method follow-
416 ing [7] to obtain a mask of the facial skin area.

417 **Implementation Details.** We implement our model in
418 PyTorch [39], using the differentiable rasterizer from Py-
419 torch3D [46] for rendering. We employ the Adam [29]
420 optimizer with a learning rate of $1e-4$ for the base model
421 and $1e-3$ for the person-specific transfer. Input images are
422 cropped and aligned with RetinaFace [13], and resized to
423 256×256 . Additional details on data augmentation, hyper-
424 parameter settings (e.g., loss balancing weights), and expla-
425 nations of the losses are in the supplementary materials.

426 4.2. Quantitative Comparison

427 We compare the accuracy of our models in face align-
428 ment with publicly available facial reconstruction methods,
429 namely 3DDFA-v2 [20], SynergyNet [76], DECA [18] and
430 EMOCA [12]. To comprehensively demonstrate the supe-
431 riority of our method in coarse shape, encompassing fa-
432 cial contours and features, we conduct evaluations across
433 monocular image reconstruction (300-W dataset [55]),
434 monocular video reconstruction (300-VW dataset [61]), and
435 multi-view image reconstruction (FaceScape dataset [78]).
436 Note that we do not evaluate our method on 3D bench-
437 marks, as mainstream self-supervised face reconstruc-
438 tion approaches typically assume an orthographic camera
439 model, whereas 3D dataset photos are often taken from
440 close distances and exhibit significant perspective distor-
441 tion. Thus, directly comparing orthographic projection-
442 based reconstructions with ground truth would be unfair.

443 **300-W Dataset.** We employed the 300-W dataset to as-
444 sess the precision of our base model in single-image face
445 alignment. As shown in Tab.1 and Fig. 5, on 1424 cleaned
446 test images, our method achieves a lower RMSE error [55]
447 than previous works for both boundary and inner land-
448 marks. This is attributed to our novel silhouette vertex re-
449 projection loss, which establishes more precise correspon-
450 dences for the ground-truth 2D silhouette landmarks while
451 naturally mitigating the relatively large variance in manual
452 landmark annotations along the silhouette tangent.

453 **300-VW Dataset.** The 300-VW dataset provides a com-
454 prehensive benchmark for landmark tracking in long-term
455 ‘in-the-wild’ facial videos. Due to the semi-automatic an-
456 notation utilized in the 300-VW Challenge [9], discrepan-
457 cies exist between the annotated landmarks and their ac-
458 tual facial positions. Consequently, testing on the 300-VW
459 dataset serves primarily as a reference for evaluating face
460 reconstruction accuracy in videos with continuously chang-
461 ing poses. Methods exhibiting similar test errors should
462 be considered comparable. As depicted in Tab.1, our base
463 model delivers results on par with DECA, EMOCA-v2, and
464 3DDFA-v2 [20], while significantly surpassing Synergy-
465 Net [76]. Additionally, our person-specific models substan-

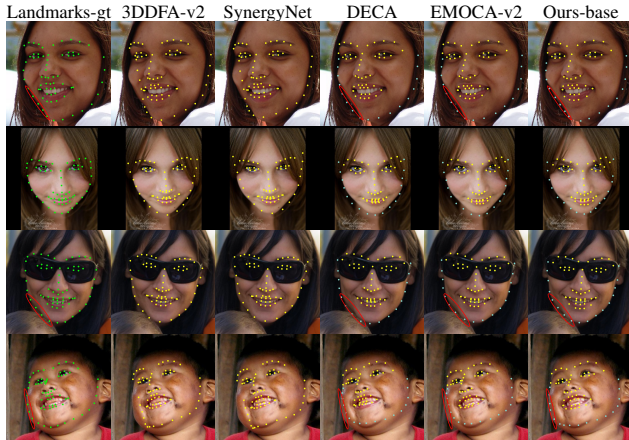


Figure 5. **Face alignment on 300-W [55]**. From left to right: Ground truth 2D landmarks, projected 3D landmarks estimated by 3DDFA-v2 [20] and SynergyNet [76], and '2D landmarks' provided by DECA [18], EMOCA-v2 [12], and our base model.

tially outperform both our base model and previous works.

FaceScape Dataset. FaceScape [78] is a large-scale detailed 3D face dataset with multi-view images, camera parameters, 3D face scans, and parametric models with their registration parameters. To address the significant perspective distortion in the images, we used FaceScape’s parametric models to extract 3D facial landmarks and projected them onto the image plane using the camera parameters, establishing them as the ground truth for image facial landmarks. As only 43 landmarks from FaceScape’s parametric model are applicable (numbered 17 ~ 59 in the 68-landmark annotation), we completed the set (17 facial boundary and 8 inner mouth circle landmarks) with the annotation from HRNet [64]. As shown in Tab.1, our base model outperforms previous works in both inner and overall landmark accuracy. The transferred person-specific models further reduce the error across different views for each identity significantly.

4.3. Qualitative Comparison

Given an in-the-wild image, our base model reconstructs a 3D face with animatable details. Given multiple images or a video of an individual, we transfer the base model by inserting trainable person-specific adapters. Our person-specific model achieves higher fidelity reconstruction of images from that individual. We conduct visualized comparisons with previous work in terms of self-supervised coarse shape reconstruction [12, 18, 20, 76], detail reconstruction [12, 18, 72, 78], and detail animation [12, 18]. The input images are taken from the FaceForensics++ dataset [52], where the images and identities were never encountered during the training of the base model.

Coarse Shape Reconstruction. Fig. 6 qualitatively compares the results of our base and person-specific models with state-of-the-art coarse reconstruction methods [12, 18,

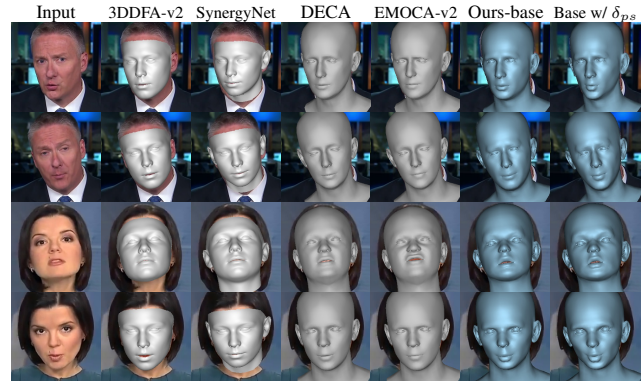


Figure 6. **Comparison on coarse shape reconstruction.** From left to right: Input image, 3DDFA-v2 [20], SynergyNet [76], DECA [18], EMOCA-v2 [12], our base model, and our transferred person-specific models.

20, 76] that are publicly available. Compared to these methods, our base model exhibits higher accuracy in fitting the outer contour, pose, and facial feature representation. Our person-specific model further enhances these advantages.

Detailed Reconstruction. Fig. 7 visually compares our work to existing detailed reconstruction methods [12, 18, 72, 78]. Several methods [72, 78] optimize for the current image, which limits inference speed and robustness to pose and occlusion. Previous generic models that reconstruct animatable geometric details [12, 18] struggle with the fidelity of person-specific facial details, as seen in Fig. 7. Compared to the base model (penultimate column), the transferred person-specific models (last column) exhibit improved accuracy in wrinkle details.

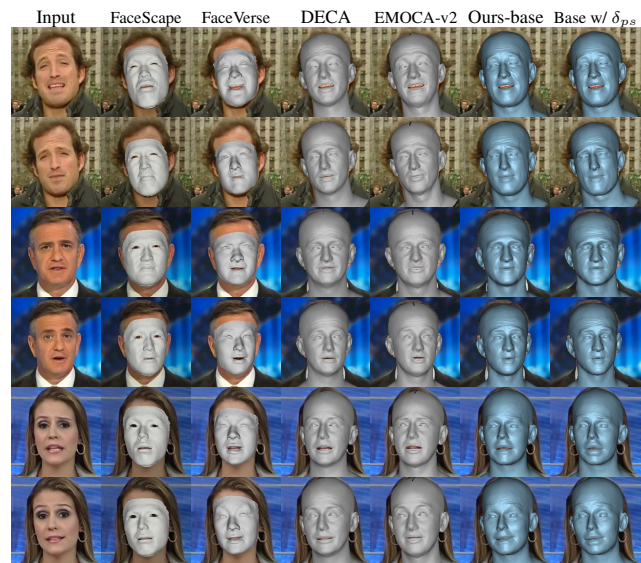


Figure 7. **Comparison on detail shape reconstruction.** From left to right: Input image, FaceScape [78], FaceVerse [72], DECA [18], EMOCA-v2 [12], our base model, and our person-specific models.

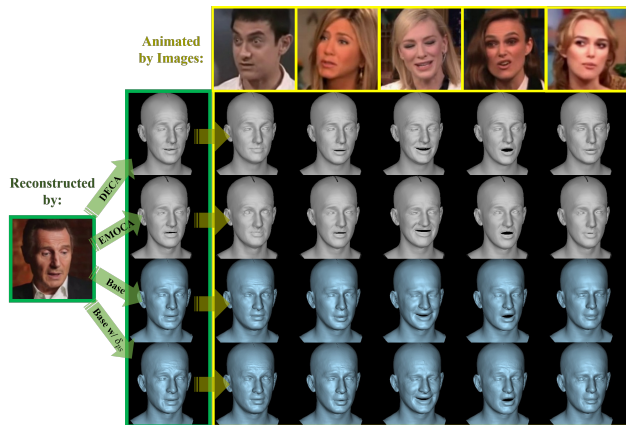


Figure 8. **Comparison on face animation.** Given a source image, DECA [18] (row 2), EMOCA-v2 [12] (row 3), and our base (row 4) and person-specific (row 5) models can respectively generate detailed 3D faces (green boxes). With a driving image (yellow boxes), these models can drive the face to exhibit corresponding expressions.

514 **Detailed Face Animation.** Fig. 8 demonstrates the anima-
 515 tion quality of our models compared to the state-of-the-
 516 art detail animation models [12, 18]. Our base model has
 517 learned a rich prior of expression-related details, surpass-
 518 ing existing works in realism and accuracy. Meanwhile,
 519 our transferred model captures person-specific details more
 520 finely while inheriting the animatability and robustness to
 521 in-the-wild driving images from the base model. This en-
 522 hances the model’s accuracy in reconstructing the specific
 523 individual, with more detailed and enriched features.

524 4.4. Ablation Studies

525 **Silhouette Vertex Re-Projection Loss.** We trained a
 526 network, Ours-base (w/o \mathcal{L}_{sil}), without the proposed sil-
 527 houette vertex re-projection loss \mathcal{L}_{sil} , using FLAME’s
 528 landmark marching algorithm [33] to apply the landmark
 529 re-projection loss across all 68 landmarks, as done by
 530 DECA [18]. The penultimate row in Table 1 shows the eval-
 531 uation results of Ours-base (w/o \mathcal{L}_{sil}) on the 300-W [55]
 532 and 300-VW [61] dataset. Ours-base (w/o \mathcal{L}_{sil}) performs
 533 slightly worse than DECA on the 300-W dataset in terms
 534 of overall error, which might be attributed to the different
 535 choice of training data (DECA uses VGGFace2 [4] and
 536 VoxCeleb2 [10]). In contrast, using \mathcal{L}_{sil} improves edge and
 537 interior landmark errors by 9.1% and 13.2%, respectively,
 538 on 300-W. Fig. 9 visually shows the contribution of \mathcal{L}_{sil} , in
 539 terms of boundary fitting accuracy.

540 **Teacher-Student Loss.** We present an ablation study
 541 on the proposed teacher-student strategy for training the de-
 542 tail network. Fig. 9 demonstrates the contribution of the
 543 teacher-student strategy to the facial detail reconstruction.
 544 The network trained without the teacher supervision loss
 545 \mathcal{L}_{Tchr} (Equation 5) (with other settings unchanged) gener-

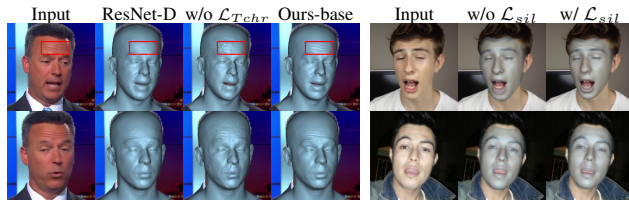


Figure 9. **Ablation studies.** Left: Compared to MAE (Ours-base), using convolutional (ResNet) and deconvolutional networks (ResNet-D) struggles to capture expression-dependent details. When training MAE without incorporating the teacher supervision loss \mathcal{L}_{Tchr} (Eqn. 5), it results in inaccurate wrinkles and artifacts. Right: Without \mathcal{L}_{sil} , the facial boundary does not properly align with the input image.

546 ates facial details with numerous unrealistic artifacts. This
 547 occurs because the shape and the rendered RGB image do
 548 not have a one-to-one correspondence, resulting in fewer
 549 constraints on the optimization direction when training the
 550 network using shape-from-shading loss, making it challeng-
 551 ing to achieve acceptable results.

552 **Network Architecture.** We compared the effectiveness
 553 of using convolutional (ResNet) and deconvolutional net-
 554 works versus a MAE for detail reconstruction, both employ-
 555 ing the teacher-student loss. For the former, we adopted
 556 the same network architecture and dynamic detail driving
 557 method as DECA [18]. Fig. 9 demonstrates that the MAE
 558 captures expression-related details more effectively. This is
 559 due to the superior long-range dependency capture and fea-
 560 ture extraction capabilities of the MAE architecture we em-
 561 ployed. Additionally, the transformer structure allows us to
 562 insert adapter layers [23], enabling an incremental person-
 563 specific transfer to retain the generalization capabilities of
 564 the base model on face animation and occlusion.

565 5. Conclusion

566 We propose constructing person-specific 3D face recon-
 567 struction models by integrating lightweight adapters into a
 568 large-scale ViT-MAE base model. During the coarse recon-
 569 struction stage, a novel silhouette vertex re-projection loss
 570 is introduced to address the issue of “landmark marching”,
 571 thereby correcting the misalignment of boundary facial
 572 landmarks and achieving state-of-the-art performance. In
 573 the detailed stage, a teacher-student loss is employed to
 574 resolve the ambiguities inherent in the self-supervised
 575 shape-from-shading approach, allowing the detail MAE
 576 to effectively capture rich and accurate features. When
 577 provided with multiple images or videos of an individual,
 578 we can further transfer the base model to a person-specific
 579 model, improving reconstruction accuracy and enabling
 580 more effective decoupling of identity and expression
 581 details. Our advantages in facial boundary and detail align-
 582 ment, combined with the ability to animate details through
 583 facial movements, make our approach highly suitable for
 584 face animation, wrinkle transfer, and downstream applica-
 585 tions such as face reenactment and virtual avatar creation.

586

587

References

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20311–20322, 2022. 1, 4, 6, 7, 8

642

643

644

[13] Jiankang Deng, Jia Guo, Yuxiang Zhou, Jinke Yu, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-stage dense face localisation in the wild. *arXiv preprint arXiv:1905.00641*, 2019. 6

645

646

647

648

[14] Qixin Deng, Luming Ma, Aobo Jin, Huikun Bi, Binh Huy Le, and Zhigang Deng. Plausible 3d face wrinkle generation using variational autoencoders. *IEEE Transactions on Visualization and Computer Graphics*, 28(9):3113–3125, 2021. 3

649

650

651

652

653

[15] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 285–295, 2019. 1, 3, 4

654

655

656

657

658

659

[16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 4

660

661

662

663

664

665

[17] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18770–18780, 2022. 1

666

667

668

669

670

[18] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d face model from in-the-wild images. *ACM Transactions on Graphics (TOG)*, 40(4): 88:1–88:13, 2021. 1, 2, 3, 4, 5, 6, 7, 8

671

672

673

674

[19] Panagiotis P Filintisis, George Retsinas, Foivos Paraperas-Papantoniou, Athanasios Katsamanis, Anastasios Roussos, and Petros Maragos. Visual speech-aware perceptual 3d facial expression reconstruction from videos. *arXiv preprint arXiv:2207.11094*, 2022. 4, 6

675

676

677

678

679

[20] Jianzhu Guo, Xiangyu Zhu, Yang Yang, Fan Yang, Zhen Lei, and Stan Z Li. Towards fast, accurate and stable 3d dense face alignment. In *Proceedings of the European Conference on Computer Vision*, pages 152–168. Springer, 2020. 1, 2, 6, 7

680

681

682

683

684

[21] Yudong Guo, Jianfei Cai, Boyi Jiang, Jianmin Zheng, et al. Cnn-based real-time dense face reconstruction with inverse-rendered photo-realistic face images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6):1294–1307, 2018. 1, 3

685

686

687

688

689

[22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 1, 2, 4

690

691

692

693

694

[23] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *Proceedings of the International Confer-*

695

696

697

698

- 699 *ence on Machine Learning*, pages 2790–2799. PMLR, 2019.
700 2, 5, 8
- 701 [24] Liwen Hu, Shunsuke Saito, Lingyu Wei, Koki Nagano, Jae-
702 woo Seo, Jens Fursund, Iman Sadeghi, Carrie Sun, Yen-
703 Chun Chen, and Hao Li. Avatar digitization from a single
704 image for real-time rendering. *ACM Transactions on Graph-
705 ics (ToG)*, 36(6):1–14, 2017. 1
- 706 [25] Yangyu Huang, Hao Yang, Chong Li, Jongyoo Kim, and
707 Fangyun Wei. Adnet: Leveraging error-bias towards normal
708 direction in face alignment. In *Proceedings of the IEEE/CVF
709 International Conference on Computer Vision*, pages 3080–
710 3090, 2021. 2
- 711 [26] Luo Jiang, Juyong Zhang, Bailin Deng, Hao Li, and Ligang
712 Liu. 3d face reconstruction with geometry details from a
713 single image. *IEEE Transactions on Image Processing*, 27
714 (10):4756–4770, 2018. 1, 2, 3
- 715 [27] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng
716 Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian
717 Richardt, Michael Zollhöfer, and Christian Theobalt. Deep
718 video portraits. *ACM transactions on graphics (TOG)*, 37(4):
719 1–14, 2018. 1
- 720 [28] Hyeonwoo Kim, Michael Zollhöfer, Ayush Tewari, Justus
721 Thies, Christian Richardt, and Christian Theobalt. Inverse-
722 facenet: Deep monocular inverse face rendering. In *Proceed-
723 ings of the IEEE Conference on Computer Vision and Pattern
724 Recognition*, pages 4625–4634, 2018. 1
- 725 [29] Diederik P Kingma and Jimmy Ba. Adam: A method for
726 stochastic optimization. *3rd International Conference on
727 Learning Representations, ICLR 2015, San Diego, CA, USA,
728 May 7-9, 2015, Conference Track Proceedings*, 2014. 6
- 729 [30] Alexandros Lattas, Stylianos Moschoglou, Baris Gecer,
730 Stylianos Ploumpis, Vasileios Triantafyllou, Abhijeet
731 Ghosh, and Stefanos Zafeiriou. Avatarme: Realistically ren-
732 derable 3d facial reconstruction” in-the-wild”. In *Proceed-
733 ings of the IEEE/CVF conference on computer vision and
734 pattern recognition*, pages 760–769, 2020. 1
- 735 [31] Jiaman Li, Zhengfei Kuang, Yajie Zhao, Mingming He, Karl
736 Bladin, and Hao Li. Dynamic facial asset and rig generation
737 from a single scan. *ACM Transactions on Graphics (TOG)*,
738 39(6):215:1–215:18, 2020. 3
- 739 [32] Ruilong Li, Karl Bladin, Yajie Zhao, Chinmay Chinara,
740 Owen Ingraham, Pengda Xiang, Xinglei Ren, Pratusha
741 Prasad, Bipin Kishore, Jun Xing, et al. Learning forma-
742 tion of physically-based face attributes. In *Proceedings of
743 the IEEE/CVF Conference on Computer Vision and Pattern
744 Recognition*, pages 3407–3416, 2020. 3
- 745 [33] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier
746 Romero. Learning a model of facial shape and expression
747 from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 3,
748 4, 8
- 749 [34] Yue Li, Liqian Ma, Haoqiang Fan, and Kenny Mitchell.
750 Feature-preserving detailed 3d face reconstruction from a
751 single image. In *Proceedings of the 15th ACM SIGGRAPH
752 European Conference on Visual Media Production*, pages
753 1:1–1:9, 2018. 1, 3
- 754 [35] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei
755 Lyu. Celeb-df: A large-scale challenging dataset for deep-
756 fake forensics. In *Proceedings of the IEEE/CVF conference
on computer vision and pattern recognition*, pages 3207–
3216, 2020. 5, 6
- [36] Jingwang Ling, Zhibo Wang, Ming Lu, Quan Wang, Chen
Qian, and Feng Xu. Structure-aware editable morphable
model for 3d facial detail animation and manipulation. In
*Proceedings of the European Conference on Computer Vi-
sion*, pages 249–267. Springer, 2022. 3
- [37] Peng Liu, Yao Yu, Yu Zhou, and Sidan Du. Single view
3d face reconstruction with landmark updating. In *IEEE
Conference on Multimedia Information Processing and Re-
trieval*, pages 403–408. IEEE, 2019. 2
- [38] Luming Ma and Zhigang Deng. Real-time hierarchical fac-
ial performance capture. In *Proceedings of the ACM SIG-
GRAPH Symposium on Interactive 3D Graphics and Games*,
pages 11:1–11:10, 2019. 1, 3
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer,
James Bradbury, Gregory Chanan, Trevor Killeen, Zeming
Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An im-
perative style, high-performance deep learning library. *Ad-
vances in neural information processing systems*, 32, 2019.
6
- [40] 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 sig-
nal based surveillance*, pages 296–301. Ieee, 2009. 3
- [41] Marcel Pietraschke and Volker Blanz. Automated 3d face
reconstruction from multiple images using quality measures.
In *Proceedings of the IEEE Conference on Computer Vision
and Pattern Recognition*, pages 3418–3427, 2016. 3
- [42] Stylianos Ploumpis, Evangelos Ververas, Eimear O’Sullivan,
Stylianos Moschoglou, Haoyang Wang, Nick Pears,
William AP Smith, Baris Gecer, and Stefanos Zafeiriou. To-
wards a complete 3d morphable model of the human head.
*IEEE Transactions on Pattern Analysis and Machine Intelli-
gence*, 43(11):4142–4160, 2020. 1
- [43] Ravi Ramamoorthi and Pat Hanrahan. An efficient represen-
tation for irradiance environment maps. In *Proceedings of
the 28th annual conference on Computer graphics and inter-
active techniques*, pages 497–500, 2001. 3
- [44] Chirag Raman, Charlie Hewitt, Erroll Wood, and Tadas Bal-
trušaitis. Mesh-tension driven expression-based wrinkles
for synthetic faces. In *Proceedings of the IEEE/CVF Win-
ter Conference on Applications of Computer Vision*, pages
3515–3525, 2023. 2, 3, 5
- [45] Eduard Ramon, Janna Escur, and Xavier Giró-i Nieto. Multi-
view 3d face reconstruction in the wild using siamese net-
works. In *Proceedings of the IEEE/CVF International Con-
ference on Computer Vision Workshops*, pages 3096–3100,
2019. 3
- [46] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Tay-
lor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia
Gkioxari. Accelerating 3d deep learning with pytorch3d.
arXiv preprint arXiv:2007.08501, 2020. 2, 6
- [47] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan
Liu. Pirenderer: Controllable portrait image generation via
semantic neural rendering. In *Proceedings of the IEEE/CVF*
757 758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813

- 814 *International Conference on Computer Vision*, pages 13759–
815 13768, 2021. 1
- 816 [48] Alexander Richard, Michael Zollhöfer, Yandong Wen, Fer-
817 nando De la Torre, and Yaser Sheikh. Meshtalk: 3d face an-
818 imation from speech using cross-modality disentanglement.
819 In *Proceedings of the IEEE/CVF International Conference*
820 *on Computer Vision*, pages 1173–1182, 2021. 1
- 821 [49] Elad Richardson, Matan Sela, Roy Or-El, and Ron Kimmel.
822 Learning detailed face reconstruction from a single image.
823 In *Proceedings of the IEEE Conference on Computer Vision*
824 *and Pattern Recognition*, pages 5553–5562, 2017. 1, 3
- 825 [50] Sami Romdhani, Volker Blanz, and Thomas Vetter. Face
826 identification by fitting a 3d morphable model using linear
827 shape and texture error functions. In *Computer Vi-*
828 *sion—ECCV 2002: 7th European Conference on Computer*
829 *Vision Copenhagen, Denmark, May 28–31, 2002 Proceed-*
830 *ings, Part IV 7*, pages 3–19. Springer, 2002. 1
- 831 [51] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-
832 net: Convolutional networks for biomedical image segmen-
833 tation. In *Medical Image Computing and Computer-Assisted*
834 *Intervention—MICCAI 2015: 18th International Conference,*
835 *Munich, Germany, October 5–9, 2015, Proceedings, Part III*
836 *18*, pages 234–241. Springer, 2015. 2, 5
- 837 [52] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Chris-
838 tian Riess, Justus Thies, and Matthias Nießner. Faceforen-
839 sics++: Learning to detect manipulated facial images. In
840 *Proceedings of the IEEE/CVF international conference on*
841 *computer vision*, pages 1–11, 2019. 7
- 842 [53] Joseph Roth, Yiyang Tong, and Xiaoming Liu. Uncon-
843 strained 3d face reconstruction. In *Proceedings of the*
844 *IEEE conference on computer vision and pattern recogni-*
845 *tion*, pages 2606–2615, 2015. 3
- 846 [54] Joseph Roth, Yiyang Tong, and Xiaoming Liu. Adaptive
847 3d face reconstruction from unconstrained photo collections.
848 In *Proceedings of the IEEE Conference on Computer Vision*
849 *and Pattern Recognition*, pages 4197–4206, 2016. 3
- 850 [55] Christos Sagonas, Epameinondas Antonakos, Georgios Tz-
851 imiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces
852 in-the-wild challenge: database and results. *Image and Vi-*
853 *sion Computing*, 47:3–18, 2016. 2, 6, 7, 8
- 854 [56] Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael J
855 Black. Learning to regress 3d face shape and expression
856 from an image without 3d supervision. In *Proceedings of*
857 *the IEEE/CVF Conference on Computer Vision and Pattern*
858 *Recognition*, pages 7763–7772, 2019. 1, 2
- 859 [57] Nedko Savov, Minh Ngõ, Sezer Karaoglu, Hamdi Dibekli-
860 oglu, and Theo Gevers. Pose and expression robust age es-
861 timation via 3d face reconstruction from a single image. In
862 *Proceedings of the IEEE/CVF International Conference on*
863 *Computer Vision Workshops*, pages 1270–1278, 2019. 2
- 864 [58] Matan Sela, Elad Richardson, and Ron Kimmel. Unre-
865 stricted facial geometry reconstruction using image-to-image
866 translation. In *Proceedings of the IEEE International Con-*
867 *ference on Computer Vision*, pages 1585–1594, 2017. 3
- 868 [59] Soumyadip Sengupta, Angjoo Kanazawa, Carlos D Castillo,
869 and David W Jacobs. Sfsnet: Learning shape, reflectance
870 and illuminance of faces ‘in the wild’. In *Proceedings of*
the IEEE/CVF Conference on Computer Vision and Pattern
Recognition, pages 6296–6305. IEEE, 2018. 1, 3
- [60] Jiayang Shang, Tianwei Shen, Shiwei Li, Lei Zhou, Ming-
min Zhen, Tian Fang, and Long Quan. Self-supervised
monocular 3d face reconstruction by occlusion-aware multi-
view geometry consistency. In *Proceedings of the European*
Conference on Computer Vision, pages 53–70. Springer,
2020. 1, 3
- [61] Jie Shen, Stefanos Zafeiriou, Grigoris G Chrysos, Jean Kos-
saifi, Georgios Tzimiropoulos, and Maja Pantic. The first
facial landmark tracking in-the-wild challenge: Benchmark
and results. In *Proceedings of the IEEE international con-*
ference on computer vision workshops, pages 50–58, 2015.
6, 8
- [62] Zhixin Shu, Duygu Ceylan, Kalyan Sunkavalli, Eli Shecht-
man, Sunil Hadap, and Dimitris Samaras. Learning monoc-
ular face reconstruction using multi-view supervision. In
2020 15th IEEE International Conference on Automatic
Face and Gesture Recognition (FG 2020), pages 241–248.
IEEE, 2020. 1
- [63] Karen Simonyan and Andrew Zisserman. Very deep convo-
lutional networks for large-scale image recognition. *arXiv*
preprint arXiv:1409.1556, 2014. 5
- [64] Ke Sun, Yang Zhao, Borui Jiang, Tianheng Cheng, Bin Xiao,
Dong Liu, Yadong Mu, Xinggang Wang, Wenyu Liu, and
Jingdong Wang. High-resolution representations for labeling
pixels and regions. *arXiv preprint arXiv:1904.04514*, 2019.
6, 7
- [65] Supasorn Suwajanakorn, Ira Kemelmacher-Shlizerman, and
Steven M Seitz. Total moving face reconstruction. In *Pro-*
ceedings of the European Conference on Computer Vision,
pages 796–812. Springer, 2014. 3
- [66] Ayush Tewari, Michael Zollhofer, Hyeonwoo Kim, Pablo
Garrido, Florian Bernard, Patrick Perez, and Christian
Theobalt. Mofa: Model-based deep convolutional face au-
toencoder for unsupervised monocular reconstruction. In
Proceedings of the IEEE International Conference on Com-
puter Vision Workshops, pages 1274–1283, 2017. 1, 2
- [67] Ayush Tewari, Michael Zollhofer, Florian Bernard, Pablo
Garrido, Hyeonwoo Kim, Patrick Perez, and Christian
Theobalt. High-fidelity monocular face reconstruction based
on an unsupervised model-based face autoencoder. *IEEE*
Transactions on Pattern Analysis and Machine Intelligence,
42(2):357–370, 2018.
- [68] Ayush Tewari, Michael Zollhöfer, Pablo Garrido, Florian
Bernard, Hyeonwoo Kim, Patrick Pérez, and Christian
Theobalt. Self-supervised multi-level face model learning
for monocular reconstruction at over 250 hz. In *Proceed-*
ings of the IEEE Conference on Computer Vision and Pattern
Recognition, pages 2549–2559, 2018. 2
- [69] Luan Tran, Feng Liu, and Xiaoming Liu. Towards high-
fidelity nonlinear 3d face morphable model. In *Proceedings*
of the IEEE/CVF Conference on Computer Vision and Pat-
tern Recognition, pages 1126–1135, 2019. 3
- [70] Anh Tuan Tran, Tal Hassner, Iacopo Masi, and Gérard
Medioni. Regressing robust and discriminative 3d mor-
phable models with a very deep neural network. In *Proceed-*

- ings of the *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1493–1502, 2017. 1
- [71] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *Proceedings of the European Conference on Computer Vision*, pages 700–717. Springer, 2020. 5, 6
- [72] Lizhen Wang, Zhiyuan Chen, Tao Yu, Chenguang Ma, Liang Li, and Yebin Liu. Faceverse: a fine-grained and detail-controllable 3d face morphable model from a hybrid dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20333–20342, 2022. 1, 3, 7
- [73] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 692–702, 2019. 6
- [74] Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative multi-column convolutional neural networks. *Advances in neural information processing systems*, 31, 2018. 5
- [75] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5
- [76] Cho-Ying Wu, Qiangeng Xu, and Ulrich Neumann. Synergy between 3dmm and 3d landmarks for accurate 3d facial geometry. In *International Conference on 3D Vision (3DV)*, pages 453–463. IEEE, 2021. 1, 6, 7
- [77] Fanzi Wu, Linchao Bao, Yajing Chen, Yonggen Ling, Yibing Song, Songnan Li, King Ngi Ngan, and Wei Liu. Mvf-net: Multi-view 3d face morphable model regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 959–968, 2019. 3
- [78] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. Facescape: A large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 598–607, 2020. 3, 6, 7
- [79] Jae Shin Yoon, Takaaki Shiratori, Shoou-I Yu, and Hyun Soo Park. Self-supervised adaptation of high-fidelity face models for monocular performance tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4601–4609, 2019. 1, 2
- [80] Xiaoxing Zeng, Xiaojiang Peng, and Yu Qiao. Df2net: A dense-fine-finer network for detailed 3d face reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2315–2324, 2019. 1, 3
- [81] 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. 1
- [82] Xiangyu Zhu, Zhen Lei, Junjie Yan, Dong Yi, and Stan Z Li. High-fidelity pose and expression normalization for face recognition in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 787–796, 2015. 2