# SADRNet: Self-Aligned Dual Face Regression Networks for Robust 3D Dense Face Alignment and Reconstruction

Zeyu Ruan, Changqing Zou, *Member, IEEE*, Longhai Wu, Gangshan Wu, *Member, IEEE*, Limin Wang, *Member, IEEE*

*Abstract*—Three-dimensional face dense alignment and reconstruction in the wild is a challenging problem as partial facial information is commonly missing in occluded and large pose face images. Large head pose variations also increase the solution space and make the modeling more difficult. Our key idea is to model occlusion and pose to decompose this challenging task into several relatively more manageable subtasks. To this end, we propose an end-to-end framework, termed as Self-aligned Dual face Regression Network (SADRNet), which predicts a pose-dependent face and a pose-independent face. They are combined by an occlusion-aware self-alignment to generate the final 3D face. Extensive experiments on two popular benchmarks, AFLW2000-3D and Florence, demonstrate that the proposed method achieves significant superior performance over existing state-of-the-art methods.

*Index Terms*—Three-dimensional deep face reconstruction, Dense face alignment, occlusion-aware attention.

## I. INTRODUCTION

Monocular 3D face reconstruction recovers 3D facial geometry from a single-view image. Dense face alignment (e.g., [1], [2]) locates all facial vertices of a face model. They are closely related to each other, and both play important roles in broad applications such as face recognition [3], [4], normalization [5], tracking [6], swapping [7], and expression recognition [8] in computer vision and graphics.

The existing methods that simultaneously address 3D dense face alignment and face reconstruction (3D-DFAFR, for short) can be roughly grouped into two categories: model-based category and model-free category. Model-based methods infer the parameters of a parametric model [9], [10], such as a 3D morphable model (3DMM [11]), by solving a nonlinear optimization problem or directly regressing with convolutional neural networks (CNNs) [12], [13], [1], [14], [15], [16]. Recent work [17], [18], [19], [20], [18] achieves high-fidelity face shape and dense face alignment by using nonlinear 3DMM decoders to improve its representation power. Rather than using a parametric model, model-free methods obtain unrestricted 3D face structure and alignment information by directly inferring the 3D position of face vertices represented in specific forms (e.g., UV map [2], volume [21]).

Z. Ruan, G. Wu, L. Wang are with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, 210023, China (e-mail: mg1833060@smail.nju.edu.cn, gswu@nju.edu.cn, lmwang@nju.edu.cn).

C. Zou is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, 510006, China (e-mail: aaronzou1125@gmail.com).

L. Wu is with the Samsung Electronics (China) R&D Centre, Nanjing, 210012, China (e-mail: longhai.wu@samsung.com).

Although significant improvements have been achieved on the problem of 3D-DFAFR in a controlled setting during the past few years, 3D-DFAFR under unconstrained conditions is still yet to be well addressed. Specifically, under unconstrained conditions, as shown in Fig. 1, self occlusions caused by large pose orientation and inter-object occlusion of hair and glasses could significantly reduce the useful information in an image, making it challenging to generate good results. Meanwhile, large pose diversity will make it hard to distinguish the shape variations and increase the modeling difficulty. Previous works mainly tackle these problems by increasing the data size and diversity of training data [1], [2], [18], [22] or performing strong regularization on the shape [23].

To addresses the challenges of 3D-DFAFR in the wild, this paper proposes an effective solution based on the following three motivations: (1) an occluded region in the image does not contain any face information but may negatively affect the network prediction. (2) the occluded part of a face can only be inferred through the global facial structure or prior knowledge. (3) disentangling the face pose and face shape will significantly reduce the complexity of the problem of 3D-DFAFR and thus make it more tractable. Based on the above motivations, we believe *occlusion* and *pose* are two critical factors in achieving a robust 3D-DFAFR. Unfortunately, there are few 3D-DFAFR works explicitly handle with the face occlusions. The face pose estimation is also rarely discussed deeply in existing 3D-DFAFR works. Therefore, in this work, we propose to explicitly model these two factors to build a robust method for 3D-DFAFR.

We propose a self-aligned dual face regression network (SADRNet) for robust 3D dense face alignment and face reconstruction based on the above analysis. In particular, we present a dual face regression framework to decouple face pose estimation and face shape prediction in a low computational cost manner. These two regression networks share the same encoder-decoder backbone. They are equipped with its task-specific head design for predicting pose-dependent face and pose-independent face, respectively. The pose-independent shape regression could relieve the difficulty of directly predicting the original complex face shape under various poses and thus improve the shape prediction accuracy. Besides, to tackle the occlusion issue, we devise a supervised attention mechanism to enhance the discriminative features in visible areas while suppressing the occluded region's influence. The attention mechanism could be plugged into our dual face
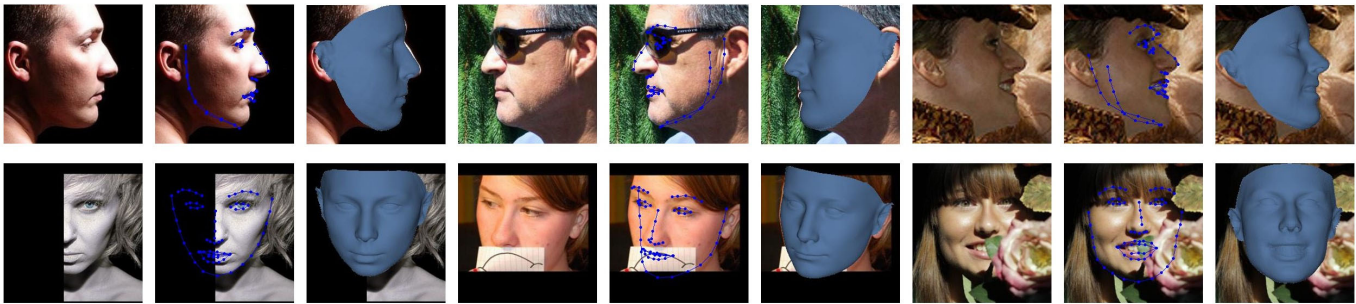
Fig. 1. Illustration of the challenges of large poses (the 1st row) and occlusions (the 2nd row). The results of both face alignment and reconstruction are demonstrated. Only 68 landmarks are plotted for better view.

regression framework to improve prediction accuracy. Finally, we propose an occlusion-aware self-alignment module to combine the pose-dependent and pose-independent faces to yield the final face reconstruction. In this alignment module, we only use the visible and sparse face landmarks to estimate the pose parameters, which could further improve the robustness of our SADRNet. Our solution significantly improves the robustness toward face occlusions in the wild. It achieves a considerable margin on both face reconstruction and dense face alignment. In summary, the main contributions of this paper are:

- We propose a self-aligned dual faces regression framework, which is robust to face pose variation and occlusion, for the problem of 3D-DFAFR.
- We propose an attention-aware mechanism for visible face region regression, which can improve the regression accuracy and robustness under various situations of face occlusion.
- The proposed end-to-end architecture is efficient and it can run at $224 \times 224/70$ FPS on a single GTX 1080 Ti GPU.
- The proposed method achieves a considerable margin on the challenging AFLW2000-3D [1] dataset and Florence 3D Faces [24] over the state-of-the-art methods.

## II. RELATED WORK

### A. 3D Face Reconstruction

Blanz and Vetter [11] proposed 3DMM to represent a face model by a linear combination of orthogonal bases obtained by PCA. In this way, 3D face reconstruction can be formulated as 3DMM parameters regression problems. Later, Paysan et al. [25] extended the model by adding more scans and decomposing the expression bases from shape bases. A lot of earlier methods regressed the 3DMM parameters by solving a nonlinear optimization function [26], [5], [25], [27], [28]. Thanks to the development of deep learning, some methods [29], [1], [30], [31], [12], [16], [32], [33] started to use deep convolutional neural network (DCNN) architectures to learn 3DMM parameters and largely replaced traditional optimization-based methods with more accurate results and shorter running time. The self-supervised training of DCNN-based methods was implemented by exploiting a differentiable renderer [34], [35], [36], [37], [38], [39], [40], which alleviated

the lack of 3D-supervised data and improved the generalization of networks. However, these model-based methods' reconstruction geometry is constrained by the linear bases with limited representation power.

Some works proposed to break the limitation by using nonlinear models [41], [42], [17], [20], [43], [41], [44]. In [19], a DCNN was used as a nonlinear 3DMM decoder. Ranjan et al. [45] used spectral graph convolutions to learn 3D faces. Zhou et al. [18] presented a nonlinear 3DMM using colored mesh decoding. Guo et al. [46] learn a more powerful nonlinear 3DMM from different data sources: scanned 3D face, RGB-D images, and RGB images.

Some other works directly obtained the full 3D geometry to avoid the restriction of parametric models and difficulty in pose estimation [47], [2], [45]. Jackson et al. [21] proposed to use a volumetric representation of 3D face shape instead of the previously used point cloud or mesh and directly regressed the voxels. Feng et al. [2] mapped the mesh of a face geometry into UV position maps and then trained a light-weighted network that obtains the 3D facial geometry along with its correspondence information.

Some works also combine the 3DMM-based regression and direct 3D geometry regression to improve the reconstruction performance. Chen et al. [48] used a 3DMM-based coarse model and a displacement map in UV space to represent a 3D face, and utilize the input image as supervision to effectively learn the facial details. Huber et al. [12] proposed to estimate an intermediate volumetric geometry and finetune it with 3DMM parameters' regression.

### B. Face Alignment

In the beginning, face alignment works aimed to locate a set of 2D facial landmarks in the image plane. Traditional works were mainly based on Active Appearance Models (AMM) and Active Shape Models (ASM) [49], [50], [51], [52], [53] and considered face alignment as a model parameter optimization problem. Then cascaded regression methods [54], [55] became popular. They iteratively refined the predictions and reached higher accuracy. With the development of deep learning, CNNs were wildly used to regress the landmarks' positions directly or predict heat maps and largely improved the performance [56], [57], [58], [59]. Since the 2D face

alignment methods have limitations on detecting invisible landmarks, the 3D face alignment problem has been widely researched in recent years. There are two major strategies for 3D face alignment: (1) separately detecting 2D landmarks and their depth [60], [61], [?], and (2) fitting a certain 3D face model [63], [64] to obtain the full 3D structure to guide the 3D landmark localization.

Since the methods mentioned above can handle only a limited quantity of landmarks, which is far from enough in some applications, 3D dense face alignment [65] started to be researched. It requires methods to offer pixel-wise facial region correspondence between two face images. As the prediction target changes from a sparse set of facial landmarks to a dense set of tens of thousands of points, it is natural to solve this problem by fitting a registered 3D face model.

*C. 3D Dense Face Alignment and Face Reconstruction*

Most model-based reconstruction approaches can be applied to dense face alignment if the face model is well registered and provides a dense correspondence [65], [17], [12]. Explicit facial pose estimation is needed in these works. Zhu et al. proposed 3DDFA [1], which is a representative work in this task. They used a cascaded CNN framework to regress the 3DMM parameters, including the pose parameters. The rotation angles are represented by a 4D unit quaternion for less difficulty in learning. However, their output faces are sensitive to the fluctuation of every parameter and hard to reach high precision. In [12], an ICP post-processing is incorporated to refine the regressed pose parameters, but it takes enormous extra computation. Zhou et al. learned a nonlinear 3DMM by directly using graph convolutions on face meshes and reached an extremely fast decoding speed. Nevertheless, the face pose is still obtained by direct regression like 3DDFA.

Some model-free methods that directly predict 3D coordinates of facial points are also applicable to the task of dense face alignment with their face representation registered with fixed semantic meaning [2], [21]. However, it is difficult for these model-free methods to handle severe occlusions and large poses since no prior knowledge or constraint is provided.

Our method's final output face geometry is represented by a UV position map, which possesses a dense correspondence between the face shape and the input image. Unlike the methods mentioned above, we explicitly deal with object occlusions by leveraging an attention mechanism on network features. We predict a pose-independent face model to avoid the large variance in face shape brought by large poses. Instead of directly regressing the pose parameters, we perform a visibility-aware self-alignment between a pose-dependent face model and a pose-independent face model to estimate the pose of a nonlinear 3DMM. The alignment process is stable as the error caused by a single outlier is apportioned by all of the landmarks. In contrast, the error of every pose parameter is accumulated in parameter-regressing methods. Besides, the alignment process is based on two regressed faces with similar shapes, rather than the matching of landmarks and a fixed face template. Therefore, the change of face shape is taken into account in pose estimation. In this way, our method is more robust and accurate than previous works.

## III. METHOD

In this section, we detail the SADRNet. We first introduce the dual face representation used in this work and the network architecture overview. We then present the occlusion-aware attention mechanism and the face fusion module. After that, we introduce the loss functions and implementation details.

*A. Dual Face Regression Framework*

**Facial geometry representation.** We assume the projection from the 3D face geometry to the 2D image is a weak perspective projection:

$$\mathbf{V} = \mathbf{Pr} * \mathbf{G} \qquad (1)$$

where $\mathbf{V}$ is the projected geometry on the 2D plane, $\mathbf{Pr} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$ is the projection matrix, $\mathbf{G} \in \mathbb{R}^{3 \times n}$ is the 3D mesh of a specific face with $n$ vertices.

We separate the 3D face geometry into pose, mean shape, and deformation as:

$$\mathbf{G} = f * \mathbf{R} * \mathbf{S} + \mathbf{t}, \qquad (2)$$

$$\mathbf{S} = (\bar{\mathbf{S}} + \mathbf{D}), \qquad (3)$$

$\mathbf{S} \in \mathbb{R}^{3 \times n}$ represents the pose-independent (i.e., pose-normalized) face shape, $\bar{\mathbf{S}} \in \mathbb{R}^{3 \times n}$ is the mean shape template provided by [25] and $\mathbf{D} \in \mathbb{R}^{3 \times n}$ is the deformation between $\mathbf{S}$ and $\bar{\mathbf{S}}$. The pose parameters consist of the scale factor $f$, the 3D rotation matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ and the 3D translation $\mathbf{t} \in \mathbb{R}^3$.

**Dual face regression and analysis.** We propose to jointly regress two face models: a pose-independent face (the face shape) $\mathbf{S}$ and a pose-dependent face $\mathbf{P}$. And then, we use a self-alignment post-process $\phi$ to estimate face pose from $\mathbf{P}$, $\mathbf{S}$ and facial vertices' visibility information $\mathbf{Vis}$:

$$\phi(\mathbf{P}, \mathbf{S}, \mathbf{Vis}) = f, \mathbf{R}, \mathbf{t}. \qquad (4)$$

The visibility information $\mathbf{Vis}$ is explained in Sec.III-B. Based on the estimated face pose $\phi(\mathbf{P}, \mathbf{S}, \mathbf{Vis})$ and pose-independent face $\mathbf{S}$, we could reconstruct our final face shape $G$ via transformation defined in Eq.2.

In fact, $\mathbf{G}$ and $\mathbf{P}$ have the same physical meaning, which means the ground truth of them are the same:

$$\hat{\mathbf{P}} = \hat{\mathbf{G}}. \qquad (5)$$

The difference is that $\mathbf{P}$ is obtained by direct network regression, while $\mathbf{G}$ is obtained by applying Eq.2. To distinguish them, we term $\mathbf{P}$ as pose-dependent face and $\mathbf{G}$ as face geometry. The learning of pose-dependent face is easy to overfit to pose and under-fit to shape as the orientation variations bring much greater point-to-point distances than the shape variations, resulting in some implausible face shape under large pose cases. By contrast, the pose-independent face $\mathbf{S}$ does not change with the pose. Thus the network would focus on the shape characteristics and learn more details. $\mathbf{S}$ is disentangled into the mean face shape template $\bar{\mathbf{S}}$ and the deformation $\mathbf{D}$ between the actual shape and the mean shape. Only the zero-centered $\mathbf{D}$ needs to be predicted. It further reduces the fitting difficulty. The mean face shape also serves
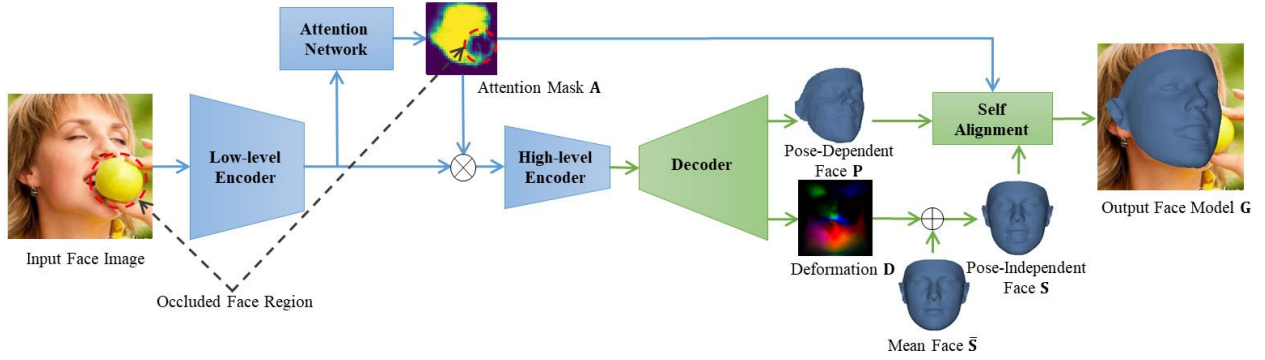
Fig. 2. The framework of our proposed self-aligned dual face regression network (SADRNet). **A** is the attention mask. **P** is the pose-dependent face. **D** is the face shape deformation (visualized in UV space). $\bar{\mathbf{S}}$ is the mean face template. **S** is the pose-independent face. **G** is the output face model.
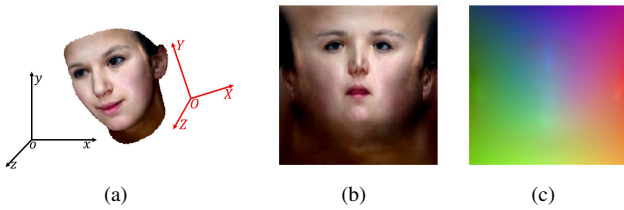


(a) (b) (c)

Fig. 3. Illustration the UV map of a 3D face model. Given a textured face model (a), we can compute the corresponding unwrapped texture in UV space (b), and the corresponding UV position map (c).
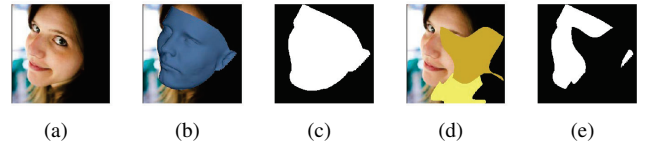


(a) (b) (c) (d) (e)

Fig. 4. An example of the synthetic attention mask annotation. (a) is the original face image. (b) is the ground truth face geometry projected on the image plane. (c) is the binary image of the projected face. (d) is the face image augmented with synthetic occlusions. (e) is the ground truth attention mask corresponding to (d).

as prior knowledge to keep the invisible facial parts plausible. By combining the shape of **S** with the pose $\phi(\mathbf{P}, \mathbf{S}, \mathbf{Vis})$ estimated from **P** and **S**, we are able to get a much better **G**, as demonstrated in experiments.

**UV map representation.** The face geometry **G**, pose-independent face **S**, mean face $\bar{\mathbf{S}}$, deformation **D**, and pose-dependent face **P** are transformed into UV space [2] as UV maps. UV map **U** is a 2D representation of 3D vertices on the face mesh model. It can be expressed as

$$\mathbf{U}(u_i, v_i) = (x_i, y_i, z_i), \quad (6)$$

where $(x_i, y_i, z_i)$ is the 3D coordinate of vertex $i$ on the 3D face mesh and $(u_i, v_i)$ is the corresponding 2D UV coordinate.

The mapping relationship between the 3D object coordinate $(X_i, Y_i, Z_i)$ and the UV coordinate $(u_i, v_i)$ of a facial vertex $i$ can be formulated as

$$u_i \rightarrow \alpha_1 \cdot Y_i + \beta_1, \quad (7)$$

$$v_i \rightarrow \alpha_2 \cdot \arctan(\frac{X_i}{Z_i}) + \beta_2, \quad (8)$$

where $\alpha_1$, $\alpha_2$, $\beta_2$, $\beta_1$ are scaling and translation constants. The mapping relationship is computed on the mean face mesh from the Basel Face Model (BFM) [25] and applied to all the face meshes. In this way, the points in the UV map are registered to the face mesh model. Fig. 3 illustrates the mapping for a better understanding. This representation guarantees the spatial consistency between face model and UV map, i.e., spatially neighboring points on the face model

are neighboring in the UV map. Since 2D UV maps can be processed by sophisticated CNNs, this ensures a great potential of the representation in applications of unconstrained situations. In the remainder of this paper, the face geometry, pose-independent face, mean face, deformation, and pose-dependent face are represented as UV maps. For clarity, the same notation is used for a thing in the two spaces (e.g., **G** is used to denote both the face geometry and the UV map of it).

**Network architecture.** Our self-aligned dual face regression network is an encoder-decoder-based architecture that regresses the deformation **D** and infers the pose parameters $f$, **R** and **t** to reconstruct the 3D face geometry from a single 2D face image. It consists of three sub-networks: encoder, attention side branch, and decoder.

The encoder network contains a sequence of residual blocks that first extract low-level features, which are then fed into the attention side branch. The attention network generates an attention mask **A** with the visible face region highlighted. Then encoder network further encodes the attended features into high-level features. The decoder then decodes them into the pose-dependent face **P** and the deformation **D**. The face shape **S** is obtained by Eq. 3. The pose parameters $f$, **R** and **t** are then obtained from **P**, **S** and **A** in the self-alignment module. The final face geometry **G** is generated by applying Eq. 2.

## B. Occlusion-Aware Attention Mechanism

To deal with the occlusions, we adopt an attention mechanism to extract features mainly from the visible face regions in the input images. This sub-network is a side branch consisting of five convolutional layers. It outputs a soft attention mask assigning high values to the pixels corresponding to the visible regions and low values to the pixels in the occluded regions and background region. The attention is applied to the low-level features to make a trade-off between resolution and accuracy.

Considering the attention mask may be inaccurate in some cases, we do not discard the information of the regions with low attention entirely. We instead highlight the features of the visible face regions according to the attention values. This operation can be formulated as

$$\mathbf{F_a} = \mathbf{F_l} \odot \exp(\mathbf{A}), \tag{9}$$

where $\mathbf{F_a}$ denotes the obtained weighted feature map, $\mathbf{F_l}$ is the low-level feature map, and $\mathbf{A}$ is the attention mask.

Once we get the attention mask, we are able to estimate the visibility of the vertices in the pose-dependent face through

$$\mathbf{Vis}(i) = \begin{cases} 0 & \text{if } n_z < 0 \\ \mathbf{A}(\lfloor x_i \rfloor, \lfloor y_i \rfloor) & \text{if } n_z \geq 0, \end{cases} \tag{10}$$

where the 3D position of vertex $i$ is $(x_i, y_i, z_i)$. The normal direction of vertex $i$ that can be obtained from its adjacent vertices is $(n_x, n_y, n_z)$. The vertex visibility is employed in the self-alignment module and is discussed detailedly in Sec. III-C.

Since no database has the ground truth for face occlusion annotation, we simulate occlusions through data augmentation. Specifically, we project the 3D face geometry to the image plane to generate a binary map that indicates the full face region. Then we overlay patterns with random shapes to the input image and use them as the occlusions as shown in Fig. 4. We use the resulting binary image as the ground truth of the attention mask. In this way, we can obtain the ground truth data for the training of the attention branch. Although occlusions in the real world can be very diverse, our experiments find the proposed method has the ability to predict real-world occlusions as shown in Fig. 5.

## C. Self-Alignment Module

This module aims to extract the pose information from the pose-dependent face $\mathbf{P}$ and the shape information from the pose-independent face $\mathbf{S}$ (see Fig. 2 for the two faces $\mathbf{P}$ and $\mathbf{S}$). It is achieved by estimating the similarity transformation matrices between $\mathbf{P}$ and $\mathbf{S}$ (i.e., $f$, $\mathbf{R}$, and $\mathbf{t}$ in Eq. 2). The similarity transformation matrix is estimated using two sets of correspondent landmarks, $\mathbf{K_S}$ and $\mathbf{K_P}$, extracted from $\mathbf{S}$ and $\mathbf{P}$, respectively. The pixels with the same coordinates in the UV maps of $\mathbf{S}$ and $\mathbf{P}$ are semantically correspondent (e.g., the nose tip corresponds to two locations with the same coordinates in the UV maps of $\mathbf{S}$ and $\mathbf{P}$). Thus, in our method, $\mathbf{K_S}$ and $\mathbf{K_P}$ are extracted based on the same UV coordinates set. Specifically, we extract the 68 landmarks as the same as [1] for $\mathbf{K_S}$ and $\mathbf{K_P}$.

We do not simply use all the 68 landmarks. We propose to use more reliable facial vertices, i.e., the visible landmarks, for the alignment. We define a diagonal matrix $\mathbf{W} \in \mathbb{R}^{k \times k}$ that represents the weight of each landmark with the visibility. The weight of the $i$th landmark is formulated as

$$\mathbf{W}(i, i) = \mathbf{Vis}(i) + eps, \tag{11}$$

where $\mathbf{Vis}(i)$ denotes the visibility of the $i$th landmark as mentioned earlier. $eps$ is set to $0.1$ in our implementation to avoid the divide-by-zero error caused by the estimated visibility of all landmarks being zero. The estimation of the similarity transformation goes as the following steps: we first estimate $f$, and then estimate $\mathbf{R}$ and $\mathbf{t}$ using the singular value decomposition method. Specifically, We first compute the two weighted centroids $\mathbf{M_S}$ and $\mathbf{M_P}$ of the two sets of landmarks $\mathbf{K_S}$ and $\mathbf{K_P}$ by

$$\mathbf{M_S} = \frac{\sum_{i=1}^{k} \mathbf{W}(i, i) * \mathbf{K_S}(i)}{|\mathbf{W}|}, \tag{12}$$

$$\mathbf{M_P} = \frac{\sum_{i=1}^{k} \mathbf{W}(i, i) * \mathbf{K_P}(i)}{|\mathbf{W}|}. \tag{13}$$

Then we obtain $f$ by

$$f = \frac{\sum_{i=1}^{k} \|\mathbf{K_P}(i) - \mathbf{M_P}\|}{\sum_{i=1}^{k} \|\mathbf{K_S}(i) - \mathbf{M_S}\|}. \tag{14}$$

After that, we normalize $\mathbf{K_S}$ and $\mathbf{K_P}$ as bellow:

$$\mathbf{K_S'} = f * (\mathbf{K_S} - \mathbf{M_S}), \tag{15}$$

$$\mathbf{K_P'} = \mathbf{K_P} - \mathbf{M_P}. \tag{16}$$

Performing SVD to $\mathbf{H} = \mathbf{K_S'} * \mathbf{W} * \mathbf{K_P'}^{\mathrm{T}}$, we have $[U, \Sigma, P] = \mathrm{SVD}(\mathbf{H})$, the rotation matrix $\mathbf{R}$ between $\mathbf{K_S'}$ and $\mathbf{K_P'}$ can be obtained by

$$\mathbf{R} = P * U^{\mathrm{T}}, \tag{17}$$

and $\mathbf{t}$ can be obtained by

$$\mathbf{t} = \mathbf{M_P} - \mathbf{R} * \mathbf{M_S}. \tag{18}$$

The pose estimation from landmarks requires only a small amount of computation and has high accuracy because the estimation is performed on a number of reliable landmarks on the two regressed faces.

## D. Loss Functions

The loss of our SADRNet consists of 6 components, the face geometry loss $\mathcal{L}_G$, the deformation loss $\mathcal{L}_D$, the pose-dependent face loss $\mathcal{L}_P$, the attention mask loss $\mathcal{L}_A$, the edge length loss $\mathcal{L}_E$, and the normal vector loss $\mathcal{L}_V$. Since the significance of the vertices in different face regions may differ considerably, we adopt a weight mask $\mathbf{M}$ as [2] to $\mathcal{L}_D$, $\mathcal{L}_P$, and $\mathcal{L}_G$. We adjust the weight ratio of the four sub-regions for better training results. We use 16:12:3:0 for sub-region1 (the landmarks): sub-region2 (the eyes, nose, and mouth): sub-region3 (the cheek, chin and forehead): sub-region4 (the neck).
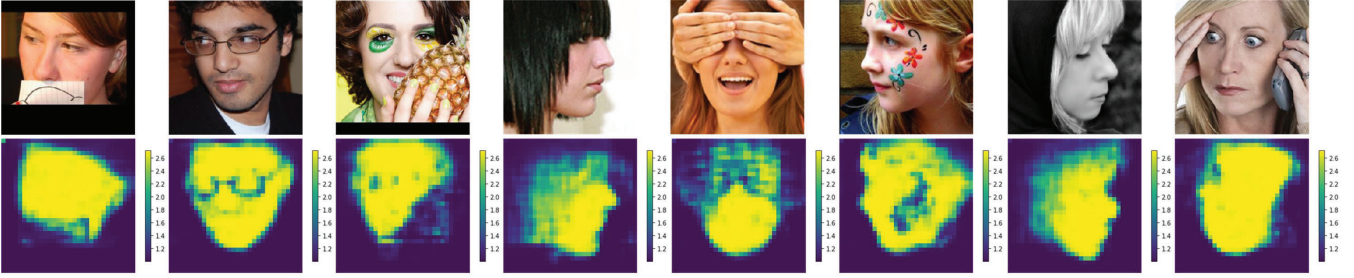
Fig. 5. The top row shows the input images and the bottom row shows the attention masks predicted by our attention branch.

$\mathcal{L}_G, \mathcal{L}_D$, and $\mathcal{L}_P$ can be obtained by the weighted average Euclidean distance between the estimated value and the ground truth, i.e., generalized as

$$\mathcal{L}_N = \sum_{u=1}^{h}\sum_{v=1}^{w}\|\mathbf{N}(u,v) - \hat{\mathbf{N}}(u,v)\|_2 \cdot \mathbf{M}(u,v). \quad (19)$$

$N$ denotes any one of the face geometry $\mathbf{G}$, the pose-dependent face $\mathbf{P}$, and the shape deformation $\mathbf{D}$. $h$ and $w$ are the height and width of the maps. $\mathbf{N}(u,v)$ denotes the estimated 3D coordinates of the vertex at the UV location $(u,v)$ in the UV maps of $\mathbf{G}$, $\mathbf{P}$, or $\mathbf{D}$. The symbol ˆ denotes the corresponding ground truth. $\mathbf{M}(u,v)$ denotes the weight value at the pixel location $(u,v)$ in $\mathbf{M}$.

We use binary cross entropy (BCE) between the predicted attention mask $\mathbf{A}$ and the ground truth $\hat{\mathbf{A}}$ to compute $\mathcal{L}_A$.

The edge length loss $\mathcal{L}_E$ is defined based on the lengths of the edges in the pose-independent face $\mathbf{S}$ as

$$\mathbf{E}_{ij} = \mathbf{S}(u_i, v_i) - \mathbf{S}(u_j, v_j), \quad (20)$$

$$\mathcal{L}_E = \sum_{(i,j)\in\mathcal{E}} |\|\mathbf{E}_{ij}\|_2 - \|\hat{\mathbf{E}}_{ij}\|_2|. \quad (21)$$

$\mathcal{E}$ is the edges set. It is composed of every pair of adjacent vertices $(i,j)$ in the UV map. $(u_i, v_i)$ is the UV coordinates of vertex $i$. $\mathbf{E}_{ij}$ is the edge vector.

The normal vector loss $\mathcal{L}_V$ is defined as

$$\mathbf{n}_{ijk} = \frac{\mathbf{E}_{ij} \times \mathbf{E}_{jk}}{\|\mathbf{E}_{ij} \times \mathbf{E}_{jk}\|_2}, \quad (22)$$

$$\begin{aligned}
\mathcal{L}_V = \sum_{(i,j,k)\in\mathcal{T}} & \left|\left\langle \frac{\mathbf{E}_{ij}}{\|\mathbf{E}_{ij}\|_2}, \mathbf{n}_{ijk}\right\rangle\right| \\
& + \left|\left\langle \frac{\mathbf{E}_{jk}}{\|\mathbf{E}_{jk}\|_2}, \mathbf{n}_{ijk}\right\rangle\right| \quad (23) \\
& + \left|\left\langle \frac{\mathbf{E}_{ki}}{\|\mathbf{E}_{ki}\|_2}, \mathbf{n}_{ijk}\right\rangle\right|,
\end{aligned}$$

where $\mathcal{T}$ is the triangle facets set, $\mathbf{n}_{i,j,k}$ is the normal vector of the triangle facet $(i,j,k)$. The edge length loss and the normal vector loss are used to generate better-looking face models.

The entire loss function of SADRNet is given by

$$\mathcal{L} = \beta_G\mathcal{L}_G + \beta_D\mathcal{L}_D + \beta_P\mathcal{L}_P + \beta_A\mathcal{L}_A + \beta_E\mathcal{L}_E + \beta_V\mathcal{L}_V, \quad (24)$$

where $\beta_G$, $\beta_D$, $\beta_P$, $\beta_A$, $\beta_E$, $\beta_V$ are respectively set to 0.1, 0.5, 1, 0.05, 1, 0.1 in our implementation.

### E. Implementation Details

Our network takes cropped-out $256 \times 256 \times 3$ images as the input, regardless of the effect of the face detector. The network starts with a single convolution layer followed by a low-level feature extractor that consists of 6 residual blocks [66] and outputs a $32 \times 32 \times 128$ feature map. The attention sub-network contains 5 convolution layers and a sigmoid activation. The high-level feature extractor contains 4 residual blocks and outputs a $8 \times 8 \times 512$ feature map. The decoder starts with 10 transpose convolution layers that up-sample the feature map to $64 \times 64 \times 64$, followed by 7 transpose convolution layers to output a UV map of the face shape with size $256 \times 256 \times 3$ and another 7 layers to output a UV map of the pose-dependent face with size $256 \times 256 \times 3$.

We follow [1], [2], [12] and use the full *300W-LP* [1] as the training set. 300W-LP contains $122,450$ face images generated from *300W* [1] by 3D rotation around the y axis and horizontal flipping. Since there is no sample with a larger than 90 degrees yaw angle in 300W-LP, we generate $5,000$ samples with yaw angles range from 90 to 105 degrees by 3D rotation to supplement the dataset. Similar to [2], we augment the training data by randomly rotating the image from -90 to 90 degrees, translating in the range of 0 to 10 percent of input size, scaling the image size from 0.9 to 1.1, and scaling the color channel separately from 0.6 to 1.4. We also generate synthetic occlusions as explained in Sec. III-B. We use Adam optimizer with a gradual warm-up strategy [67]. We start from a learning rate of 1e-5 and increase it by a constant amount at each iteration. After 4 epochs, the learning rate reaches 1e-4. Then we use an exponential scheduler that decays the learning rate by 0.85 every epoch. The batch size is set to 16. We train our network for 25 epochs.

The original annotations of 300W-LP are 3DMM parameters, so the ground truth of $\mathbf{D}$ can be generated by computing $\mathbf{D} = \sum_i \alpha_i\mathbf{A}_i$, where $(\alpha_1, \alpha_2, \dots)$ are the shape parameters, $(\mathbf{A}_1, \mathbf{A}_2, \dots)$ are the shape basis vectors (including the identity and expression basis vectors).

## IV. EXPERIMENTS

In this section, we evaluate the performance of the proposed method in terms of 3D face reconstruction, dense face alignment and head pose estimation. Our SADRNet is quantitatively compared with state-of-the-art methods including CMD (2019 [18]), SPDT (2019 [23]), PRN (2018 [2]) and

3DDFAv2 (2020 [16]). The qualitative comparison of PRN, MGCNet [40] and our method are demonstrated in Fig. 6.

### A. Evaluation Datasets

**AFLW2000-3D** [1] is an in-the-wild dataset with large variations in pose, illumination, expression, and occlusion. There are 2,000 images annotated with 68 3D landmarks and fitted 3DMM parameters to recover the ground truth face model. We evaluate both face alignment and 3D face reconstruction performance on this dataset.

**Florence** [24] is a publicly available database of 53 subjects. The ground truth annotations are meshes scanned by a structured-light system. Similar to [21], [2], [12], each face mesh is rendered in 20 poses: a pitch angle of -15, 0, 20, or 25 degrees and a yaw angle of -80, -40, 0, 40, and 80 degrees to generate the face images. We evaluate the 3D face reconstruction performance on this dataset.

### B. Face Alignment

We employ normalized mean error (NME) as the evaluation metric. It is defined as the normalized mean Euclidean distance between each pair of corresponding points in the predicted result $\mathbf{p}$ and the ground truth $\hat{\mathbf{p}}$:

$$\text{NME} = \frac{1}{N} \sum_{i=1}^{N} \frac{\|\mathbf{p_i} - \hat{\mathbf{p_i}}\|_2}{d}. \tag{25}$$

For a fair comparison, the normalization factor $d$ of the compared methods is computed in the same way. For sparse and dense alignment, the normalization factor of NME is defined as $\sqrt{h * w}$ where $h$ and $w$ are the height and width of the bounding box of all the evaluated points.

We evaluate the performance of 2D and 3D sparse alignment on the point set of 68 landmarks. We evaluate the performance of 2D and 3D dense alignment on the point set of around 45K points selected from the largest common face region of different methods as in [2]. Following the settings of previous works [2], [1], [18] on 2D sparse alignment, we also test the images with yaw angles in $[0°, 30°)$ (1,306 samples), $[30°, 60°)$ (462 samples) and $[60°, 90°]$ (232 samples) separately, as well as the balanced subsets of each angle interval.

The quantitative results on AFLW2000-3D are shown in Table I. The results of other methods are from published papers or produced by the official open-source codes. We can see our method is superior to the other methods on most metrics. Especially on 3D and dense alignment tasks, our method has taken a clear lead. SPDT [23] generates rendered large pose training samples and uses a CycleGAN [22] to transform the rendering style. The large pose samples in our training database 300W-LP are mostly obtained by 3D image rotation with large distortion. We argue that they outperform us for sparse alignment in-between 60 and 90 degrees is because of the quality gap of training data.

Visualized results of challenging samples in AFLW2000-3D are demonstrated in Fig. 6, in which various degrees of occlusion and different face orientations are evaluated. We compare our method with two representative methods: MGCNet [40]

and PRN [2]. MGCNet is a model-based method that fits the shape and the pose parameters of 3DMM by CNN. PRN is a model-free method that directly infers 3D coordinates of face mesh vertices with the UV position map. Our pose estimation method based on the alignment using visibility is more accurate and robust than previous works.

On the samples where the faces are partially occluded, as shown in Fig. 6c and Fig. 6d, the inaccuracy of the estimated poses of MGCNet leads to misalignment. And under severe occlusion, e.g., Fig. 6e, MGCNet fails to obtain a reasonable result, while our method is still able to precisely estimate the head pose. PRN jointly obtains pose and shape and works well for the visible facial region. However, PRN tends to give results with larger errors for invisible regions (see the misaligned landmarks in Fig. 6b and Fig. 6g). It is worth noting that AFLW2000-3D [1] is semi-automatically annotated and has some inaccurate annotations in some cases. In Fig. 7, we demonstrate some examples from AFLW2000-3D that our predictions have relatively larger NME but are apparently more accurate than the ground truth. This may narrow our method's superiority margin.

### C. 3D Face Reconstruction

On this task, we use NME normalized by 3D outer interocular distance as the evaluation metric. We follow the settings of [2] to evaluate our method on AFLW2000-3D [1] and Florence [24]. For AFLW2000-3D, we use the same set of points as we do for dense face alignment. The results under different yaw angles are also reported. As for Florence, we choose the face region containing 19K points. We use iterative closest point (ICP) to align the results to the corresponding ground truth as in [2]. Table II gives the quantitative comparison. Our method is robust to pose variations and achieves around 13% improvement over the state-of-the-art method on both datasets.

As our method employs a nonlinear shape deformation decoder that is more powerful than the linear bases based representation of MGCNet, our network can reconstruct more shape variations. It is worth mentioning that the output space of PRN is also nonlinear, but its large space, caused by entangled shape and pose, increases the learning difficulty for the high-frequency details. It is shown in Fig. 6 that eyes, lips, and noses of the face models reconstructed by our method have more details than PRN. Our attention-aware mechanism also contributes to the robustness toward the occlusions. In Figs. 6b-e, our reconstructed faces maintain natural appearances while those from PRN are distorted in occluded areas.

### D. Head Pose Estimation

On this task, we use MAE (mean absolute error) of the head pose parameters as the evaluation metric:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^{N} |\mathbf{p}_i - \hat{\mathbf{p}}_i|, \tag{26}$$

where $\mathbf{p}$ represents the predicted pose parameters and $\hat{\mathbf{p}}$ represents the ground truth.

In Table III, we compare our proposed self-aligned module with the SOTA head pose estimation methods. The actual

TABLE I
PERFORMANCE COMPARISON ON AFLW2000-3D ON THE TASK OF SPARSE ALIGNMENT (68 LANDMARKS) AND DENSE ALIGNMENT (45K POINTS). THE NME (%) ARE REPORTED. IMAGES WITH DIFFERENT YAW ANGLE RANGES ARE ALSO EVALUATED SEPARATELY. "BALANCED" DENOTES THE RESULTS ON THE SUBSET WITH BALANCED DISTRIBUTION OF THE YAW ANGLES. "MEAN" DENOTES THE AVERAGE RESULTS ON THE ENTIRE DATASET. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

| Method | 68 points | | | | | | 45k points | |
| | 2D | | | | | 3D | 2D | 3D |
| | 0 to 30 | 30 to 60 | 60 to 90 | Balanced | Mean | Mean | Mean | Mean |
|---|---|---|---|---|---|---|---|---|
| 3DDFA [1] | 3.78 | 4.54 | 7.93 | 5.42 | 6.03 | 7.50 | 5.06 | 6.55 |
| 3DFAN [60] | 2.77 | 3.48 | 4.61 | 3.62 | - | - | - | - |
| DeFA [65] | - | - | - | 4.50 | 4.36 | 6.23 | 4.44 | 6.04 |
| 3DSTN [68] | 3.15 | 4.33 | 5.98 | 4.49 | - | - | - | - |
| Nonlinear 3DMM [19] | - | - | - | 4.12 | - | - | - | - |
| PRN [2] | 2.75 | 3.51 | 4.61 | 3.62 | 3.26 | 4.70 | 3.17 | 4.40 |
| DAMDN [14] | 2.90 | 3.83 | 4.95 | 3.89 | - | - | - | - |
| CMD [18] | - | - | - | 3.90 | - | - | - | - |
| SPDT [23] | 3.56 | 4.06 | **4.11** | 3.88 | - | - | - | - |
| 3DDFAv2 [16] | **2.63** | 3.42 | 4.48 | 3.51 | - | - | - | 4.18 |
| **SADRNet (ours)** | 2.66 | **3.30** | 4.42 | **3.46** | **3.05** | **4.33** | **2.93** | **4.02** |



Fig. 6. The qualitative comparison on AFLW2000-3D dataset. The estimated landmarks are in blue and the ground truth landmarks are in red. NME(%) is shown at the bottom right of each result.
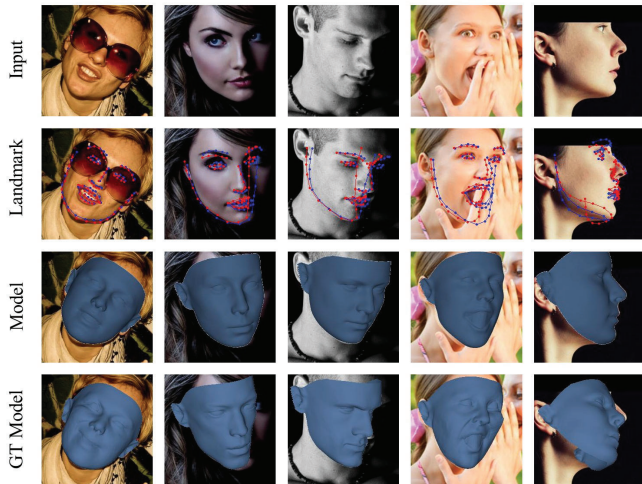
Fig. 7. Results on AFLW2000-3D from our SADRNet. Our results are more accurate than the ground truth. From the top row to the bottom row are the input images, the sparse alignment results of SADRNet and the corresponding ground truth (blue for our method and red for the ground truth), the reconstructed face models, and the ground truth face models.

TABLE II

PERFORMANCE COMPARISON ON 3D FACE RECONSTRUCTION. EVALUATIONS ARE CONDUCTED ON THE AFLW2000-3D DATASET AND FLORENCE DATASET. AROUND 45K POINTS ARE USED ON AFLW2000-3D AND 19K POINTS ARE USED ON FLORENCE. THE NME (%) NORMALIZED BY OUTER INTEROCULAR DISTANCE ARE REPORTED. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

| Method | AFLW2000-3D | | | | Florence |
|---|---|---|---|---|---|
| | 0   30 | 30   60 | 60   90 | Mean | Mean |
| 3DDFA [1] | - | - | - | 5.36 | 6.38 |
| DeFA [65] | - | - | - | 5.64 | - |
| VRN - Guided [21] | - | - | - | - | 5.26 |
| PRN [2] | 3.72 | 4.04 | 4.45 | 3.96 | 3.75 |
| SPDT [23] | - | - | - | 3.70 | 3.83 |
| 3DDFAv2 [16] | - | - | - | - | 3.56 |
| SADRNet (ours) | **3.17** | **3.42** | **3.36** | **3.25** | **3.12** |

TABLE III

PERFORMANCE COMPARISON ON HEAD POSE ESTIMATION. EVALUATIONS ARE CONDUCTED ON THE AFLW2000-3D DATASET. THE MAEs OF HEAD POSE PARAMETERS ARE REPORTED. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

| Method | AFLW2000-3D | | | |
|---|---|---|---|---|
| | yaw | pitch | roll | Mean |
| FSANet [69] | 4.50 | 6.08 | 4.64 | 5.07 |
| GLDL [70] | 3.02 | 5.06 | 3.68 | 3.92 |
| QuatNet [71] | 3.97 | 5.61 | 3.92 | 4.50 |
| FDN [72] | 3.78 | 5.61 | 3.88 | 4.42 |
| MNN [73] | 3.34 | 4.69 | 3.48 | 3.83 |
| SADRNet (ours) | 2.93 | 5.00 | 3.54 | 3.82 |
| SADRNet-fix (ours) | **2.93** | **4.43** | **2.95** | **3.44** |

TABLE IV

ABLATION STUDY ON AFLW2000-3D. THE "AT" COLUMN INDICATES USING THE ATTENTION MECHANISM OR NOT. THE "SA" COLUMN INDICATES USING THE SELF ALIGNMENT MODULE OR NOT. THE "DF" COLUMN INDICATES WHETHER 3D-DFAFR IS ACHIEVED BY REGRESSING SHAPE DEFORMATION "D" OR NOT.

| Applied approaches | | | Sparse | | Dense | | Rec. |
|---|---|---|---|---|---|---|---|
| AT | SA | DF | 2D | 3D | 2D | 3D | 3D |
| | | | 3.31 | 4.74 | 3.10 | 4.33 | 4.07 |
| | | ✓ | 3.43 | 4.94 | 3.21 | 4.50 | 3.37 |
| | ✓ | ✓ | 3.18 | 4.51 | 3.03 | 4.17 | 3.39 |
| ✓ | | | 3.16 | 4.53 | 3.03 | 4.21 | 4.02 |
| ✓ | ✓ | | 3.24 | 4.59 | 3.14 | 4.31 | 3.49 |
| ✓ | ✓ | ✓ | 3.05 | 4.33 | 2.93 | 4.02 | 3.25 |

output of our method is the transformation matrix. So we convert the rotation matrix to Euler angles to calculate MAE with the ground truth. However, when the yaw angle is close to 90 degrees, the angle conversion suffers a serious gimbal lock problem. Some examples are shown in Fig.8. Their yaw angles' errors are low, and the faces' orientations seem nice, but the errors of the pitch and roll angles are extremely high. It will make the quantitative evaluation result worse than the actual performance. Thus, besides the initial evaluation results, we also report a result marked as SADRNet-fix that ignores the samples with more than 20 degrees pitch or roll error and less than 5 degrees yaw error. However, the negatively biased performance of our method (SADRNet) is still equivalent to the best standalone head pose estimation method.
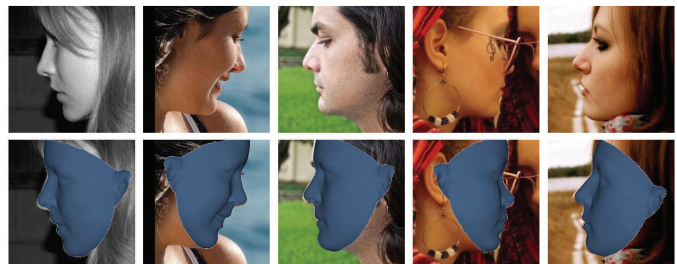


Fig. 8. Some samples in AFLW2000-3D, in which our estimated pitch angle and roll angle have an error of more than 20 degrees, while the yaw angle errors are less than 5 degrees.

### E. Ablation Study

We conduct ablation experiments on AFLW2000-3D to analyze the effectiveness of the following three modules: 1) the attention mechanism, 2) the self-alignment module, and 3) the regression via the shape deformation. The experiments are conducted on three tasks: sparse alignment, dense alignment, and 3D face reconstruction. The results are summarized in Table IV.

**Attention mechanism.** To study the contribution of the proposed attention mechanism, we remove the attention side branch from the proposed SADRNet (denoted as SADRNet-D; see Fig. 9c for the structure.) and compare it with SADRNet. The results of the two networks are respectively summarized in the 3rd and the bottom rows in Table IV. We can see that the proposed attention mechanism benefits all three tasks. In Fig. 10, we visualize the comparison. SADRNet-D is more
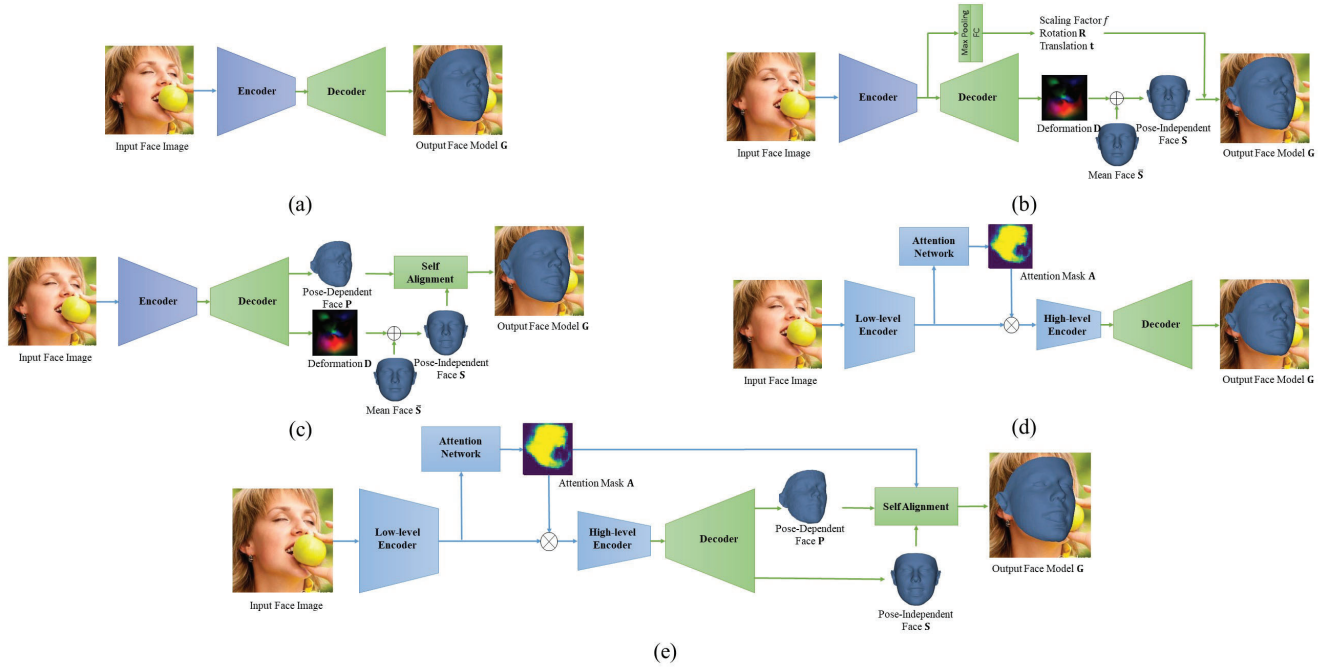
Fig. 9. Alternative frameworks of baselines in ablation study. The figures (a)-(e) correspond to the 1st-5th result rows in Table IV, respectively.

easily affected by the occlusions and has lower reconstruction accuracy in the occlusion areas. In addition, SADRNet-D cannot estimate a face orientation as accurately as SADRNet. We can see the evidence from the bottom row of Fig. 10. We also investigate the attention mechanism on the basic encoder-decoder architecture as shown in Figs. 9 a and d. The corresponding results in Table IV confirm the conclusion we draw from the comparison between SADRNet-D and SADRNet.

**Dual face regression and self alignment.** The distinguishing features of the self-alignment module lie in two points. One is that we obtain the target 3D face geometry through a two-stage refinement process (i.e., first pose-dependent face $\mathbf{P}$ and then the final 3D face model $\mathbf{G}$ as shown in Fig. 2) rather than a direct regression as shown in Fig. 9a. Another feature is that we estimate the pose information by aligning two reconstructed faces: the pose-dependent face $\mathbf{P}$ and the pose-independent face $\mathbf{S}$, rather than direct regression using the image features as in [1] (shown in Fig. 9b). By comparing the first and the third result rows in Table IV, we can observe that the two-stage refinement process brings significant gains for all of the three tasks. The comparison between the second and third result rows demonstrates that the self-alignment-based face alignment is more reliable than that based on direct pose regression from image features. However, Fig. 9b has a slightly better 3D face reconstruction than Fig. 9c. This may be because regressing $\mathbf{S}$ and $\mathbf{P}$ together may slightly affect the estimation of $\mathbf{S}$.

**Shape deformation.** This paper regresses the pose-independent face $\mathbf{S}$ through the shape deformation, which estimates the differential 3D geometry relative to the mean face template. An alternative solution to regress $\mathbf{S}$ is to directly perform the estimation in UV space from scratch by the

decoder layers as shown in Fig. 9d. By comparing the results in the 5th and bottom result rows in Table IV, it is easy to find that the shape regression based on deformation provides more accurate results than direct regression from scratch on all three tasks. Moreover, quantitatively comparing the improvement provided by using deformation (i.e., the bottom row over the 5th row) and that provided by using attention mechanism (i.e., the bottom row over the 3rd row), we can see the contribution of the deformation regression. Its improvement is more than that made by the attention mechanism in terms of 3D face reconstruction.

**Mesh loss.** Besides the losses that directly supervise the 3D coordinates of $\mathbf{P}$, $\mathbf{S}$, and $\mathbf{G}$, inspired by [74], we also adopt losses defined on the face mesh structure, i.e., the edge length loss and the normal vector loss. They do not improve the quantitative results, but help to reconstruct the face details for better visualization. In Fig.11, we demonstrate some reconstruction results of the model trained with and without the mesh loss. In the demonstrated cases, the estimated landmarks are almost the same. However, the model with mesh loss can better reflect the difference between identities. In the first example with the mesh loss, the mouth's openness is more suitable. In the second example, the curvature of the cheek is better, and the reconstructed model is more recognizable. In the third example, the hollow is more obvious.

### F. Running Time and Model Size

Our method has a model size of 60MB (including the mean face template parameters). In Table IV-F we compare the model sizes of the proposed method and the baseline methods. Our model is lighter than all other deep-learning-based methods. The network inference for one face image takes 12ms on a GTX 1080 Ti GPU. The self-alignment post
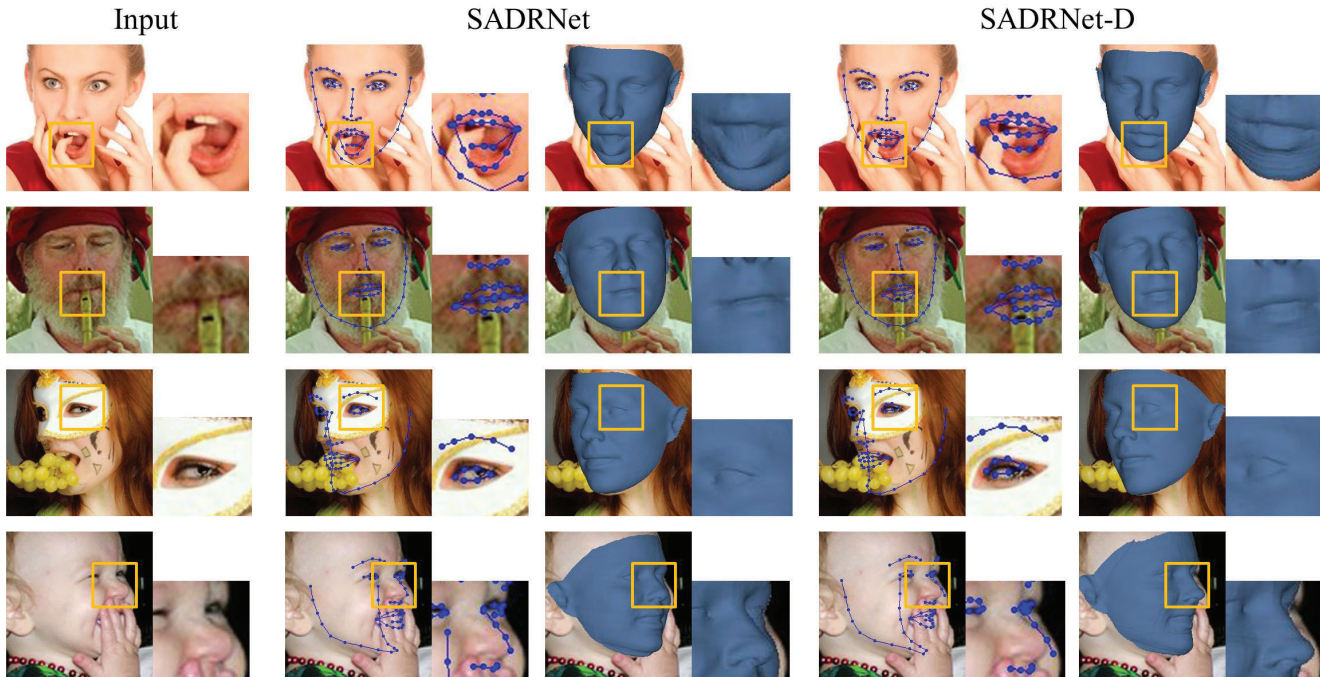
Fig. 10. SADRNet vs. SADRNet-D; results for sparse alignment and face reconstruction are demonstrated. Facial regions around occlusions are zoomed in for better visual comparison.
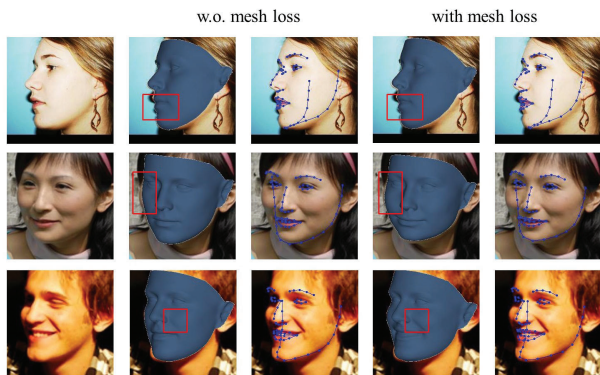


Fig. 11. With mesh loss vs. without mesh loss; results for sparse alignment and face reconstruction are demonstrated.

TABLE V
COMPARISON OF MODEL SIZE.

| Method | Size |
|---|---|
| VRN [21] | 1.5GB |
| PRN [2] | 153MB |
| Nonlinear-3DMM [17] | 152MB |
| CMD [18] | 93MB |
| SADRNet (ours) | 60MB |

process (i.e., the step generating $\mathbf{G}$ from $\mathbf{S}$ and $\mathbf{P}$) takes 4ms on GPU or 1.5ms on an Intel Xeon E5-2690 CPU @ 2.60GHz. The fastest implementation of our method reconstructs the 3D face model from a cropped image in up to 13.5 ms.

The backbone of our method, (i.e. the framework in Fig.9.a) has a model size of 52MB. The extra size introduced by the attention side branch is 7MB. The parameters of the 7-layer pose-dependent face decoder and the 7-layer pose-independent face decoder add up to 1MB.

### G. Limitations

We note the following limitations of our work:

- Our network only regresses the shape geometry and pose, but does not reconstruct the facial texture from the input image.

- The learning of the proposed SADRNet is fully supervised and depends on the costly face mesh annotation. Designing a weakly supervised architecture that can utilize additional data modalities (e.g., facial keypoint detection data, silhouette data, segmentation data) may improve the application potential.

- There is still much room for improvement in the reconstruction of high-frequency facial details (i.e. the pores, blemishes, and wrinkles) in our work. We believe the improvements can be made in two aspects. First, the training dataset we currently use is labeled with the parameters of 3DMM. The ground truth face models lack high-frequency details. Using datasets with high-fidelity ground truth models may help the detail synthesis. Considering that such training data is expensive, a self-supervised architecture that utilizes the facial details in the input image as supervision may be another feasible way. Second, although we have regressed the face shape deformation separately in the framework, the high-frequency details are relatively minor compared to the deformation and easy to ignore in learning. Further decomposing the shape and details or iteratively updating

the face shape with cascaded regressors can more adequately supervise the learning of high-frequency details. We leave these limitations as our future work.

## V. CONCLUSIONS AND FUTURE WORK

In this paper, we have proposed a self-aligned dual face regression network (SADRNet) to solve the problem of 3D face reconstruction and dense alignment under unconstrained conditions. We decouple the framework of face reconstruction to two regression modules for pose-dependent and pose-independent face shape estimation, respectively. Then, a novel self-alignment module is presented to transform the detailed and more accurate face shape into its corresponding pose view to yield the final face reconstruction. To make our method robust to occlusion, we incorporate an attention module to enhance the visible facial information and estimate the transformation matrix only with visible landmarks. We evaluate our network on the AFLW2000-3D database and the Florence database. With the power of robustness to occlusions and large pose variation, our proposed method outperforms the state-of-the-art methods by a notable margin on both face alignment and 3D reconstruction.

Lacking high-frequency facial details is the main drawback of our method and we consider to improve our method in the future from two aspects. First, we could finetune our SADRNet on a high-fidelity 3D scanned dataset or design a self-supervised learning framework to train the network with high-frequency details in the input images. Second, we plan to present a cascade reconstruction pipeline to regress our face shape in a coarse-to-fine manner and focus on more detailed face shape regression in the latter stages.

## REFERENCES

[1] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li, "Face alignment across large poses: A 3d solution," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[2] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou, "Joint 3d face reconstruction and dense alignment with position map regression network," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[3] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "Deepface: Closing the gap to human-level performance in face verification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[4] F. Liu, Q. Zhao, X. Liu, and D. Zeng, "Joint face alignment and 3d face reconstruction with application to face recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[5] X. Zhu, Z. Lei, J. Yan, D. Yi, and S. 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 (CVPR)*, 2015.

[6] X. Xiong and F. De la Torre, "Global supervised descent method," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[7] D. Chen, Q. Chen, J. Wu, X. Yu, and T. Jia, "Face swapping: Realistic image synthesis based on facial landmarks alignment," *Mathematical Problems in Engineering*, 2019.

[8] A. Mollahosseini, D. Chan, and M. H. Mahoor, "Going deeper in facial expression recognition using deep neural networks," in *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2016.

[9] T. Li, T. Bolkart, M. J. Black, H. Li, and J. Romero, "Learning a model of facial shape and expression from 4d scans," *ACM Trans. Graph.*, 2017.

[10] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "Smpl: A skinned multi-person linear model," *ACM Trans. Graph.*, 2015.

[11] V. Blanz and T. Vetter, "A morphable model for the synthesis of 3d faces," in *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, 1999.

[12] H. Yi, C. Li, Q. Cao, X. Shen, S. Li, G. Wang, and Y.-W. Tai, "Mmface: A multi-metric regression network for unconstrained face reconstruction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[13] A. Tuan Tran, T. Hassner, I. Masi, E. Paz, Y. Nirkin, and G. Medioni, "Extreme 3d face reconstruction: Seeing through occlusions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[14] L. Jiang, X.-J. Wu, and J. Kittler, "Dual attention mobdensenet(damdnet) for robust 3d face alignment," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, 2019.

[15] J. Lin, Y. Yuan, T. Shao, and K. Zhou, "Towards high-fidelity 3d face reconstruction from in-the-wild images using graph convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[16] J. Guo, X. Zhu, Y. Yang, F. Yang, Z. Lei, and S. Z. Li, "Towards fast, accurate and stable 3d dense face alignment," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.

[17] L. Tran and X. Liu, "Nonlinear 3d face morphable model," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[18] Y. Zhou, J. Deng, I. Kotsia, and S. Zafeiriou, "Dense 3d face decoding over 2500fps: Joint texture and shape convolutional mesh decoders," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[19] L. Tran and X. Liu, "On learning 3d face morphable model from in-the-wild images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.

[20] L. Tran, F. Liu, and X. Liu, "Towards high-fidelity nonlinear 3d face morphable model," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[21] A. S. Jackson, A. Bulat, V. Argyriou, and G. Tzimiropoulos, "Large pose 3d face reconstruction from a single image via direct volumetric cnn regression," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[22] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[23] J. Piao, C. Qian, and H. Li, "Semi-supervised monocular 3d face reconstruction with end-to-end shape-preserved domain transfer," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.

[24] A. D. Bagdanov, A. Del Bimbo, and I. Masi, "The florence 2d/3d hybrid face dataset," in *Proceedings of the 2011 Joint ACM Workshop on Human Gesture and Behavior Understanding*, 2011.

[25] P. Paysan, R. Knothe, B. Amberg, S. Romdhani, and T. Vetter, "A 3d face model for pose and illumination invariant face recognition," in *IEEE International Conference on Advanced Video and Signal Based Surveillance*, 2009.

[26] Y. J. Lee, S. J. Lee, K. R. Park, J. Jo, and J. Kim, "Single view-based 3d face reconstruction robust to self-occlusion," *EURASIP Journal on Advances in Signal Processing*, 2012.

[27] P. Huber, Z.-H. Feng, W. Christmas, J. Kittler, and M. Rätsch, "Fitting 3d morphable face models using local features," in *IEEE International Conference on Image Processing (ICIP)*, 2015.

[28] S. Romdhani and T. Vetter, "Estimating 3d shape and texture using pixel intensity, edges, specular highlights, texture constraints and a prior," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.

[29] A. Jourabloo and X. Liu, "Large-pose face alignment via cnn-based dense 3d model fitting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[30] M. Sela, E. Richardson, and R. Kimmel, "Unrestricted facial geometry reconstruction using image-to-image translation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[31] A. Tuan Tran, T. Hassner, I. Masi, and G. Medioni, "Regressing robust and discriminative 3d morphable models with a very deep neural network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[32] A. Chen, Z. Chen, G. Zhang, K. Mitchell, and J. Yu, "Photo-realistic facial details synthesis from single image," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.

[33] A. Lattas, S. Moschoglou, B. Gecer, S. Ploumpis, V. Triantafyllou, A. Ghosh, and S. Zafeiriou, "Avatarme: Realistically renderable 3d facial

reconstruction "in-the-wild"," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[34] K. Genova, F. Cole, A. Maschinot, A. Sarna, D. Vlasic, and W. T. Freeman, "Unsupervised training for 3d morphable model regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[35] S. Sanyal, T. Bolkart, H. Feng, and M. J. Black, "Learning to regress 3d face shape and expression from an image without 3d supervision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[36] B. Gecer, S. Ploumpis, I. Kotsia, and S. Zafeiriou, "Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[37] F. Wu, L. Bao, Y. Chen, Y. Ling, Y. Song, S. Li, K. N. Ngan, and W. Liu, "Mvf-net: Multi-view 3d face morphable model regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[38] A. Tewari, F. Bernard, P. Garrido, G. Bharaj, M. Elgharib, H.-P. Seidel, P. Perez, M. Zollhofer, and C. Theobalt, "Fml: Face model learning from videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[39] Y. Deng, J. Yang, S. Xu, D. Chen, Y. Jia, and X. 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 (CVPR) Workshops*, 2019.

[40] J. Shang, T. Shen, S. Li, L. Zhou, M. Zhen, T. Fang, and L. Quan, "Self-supervised monocular 3d face reconstruction by occlusion-aware multi-view geometry consistency," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.

[41] F. Liu, R. Zhu, D. Zeng, Q. Zhao, and X. Liu, "Disentangling features in 3d face shapes for joint face reconstruction and recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[42] T. Bagautdinov, C. Wu, J. Saragih, P. Fua, and Y. Sheikh, "Modeling facial geometry using compositional vaes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[43] Z.-H. Jiang, Q. Wu, K. Chen, and J. Zhang, "Disentangled representation learning for 3d face shape," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[44] S. Laine, T. Karras, T. Aila, A. Herva, S. Saito, R. Yu, H. Li, and J. Lehtinen, "Production-level facial performance capture using deep convolutional neural networks," in *ACM Transactions on Graphics (Proceedings of SIGGRAPH)*, 2017.

[45] A. Ranjan, T. Bolkart, S. Sanyal, and M. J. Black, "Generating 3d faces using convolutional mesh autoencoders," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[46] Y. Guo, L. Cai, and J. Zhang, "3d face from x: Learning face shape from diverse sources," *IEEE Transactions on Image Processing*, 2021.

[47] X. Fan, S. Cheng, K. Huyan, M. Hou, R. Liu, and Z. Luo, "Dual neural networks coupling data regression with explicit priors for monocular 3d face reconstruction," *IEEE Transactions on Multimedia*, 2020.

[48] Y. Chen, F. Wu, Z. Wang, Y. Song, Y. Ling, and L. Bao, "Self-supervised learning of detailed 3d face reconstruction," *IEEE Transactions on Image Processing*, 2020.

[49] a. C. T. Timothy Cootes, Gareth Edwards, "Active appearance models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2001.

[50] T. F. Cootes, C. J. Taylor, D. H. Cooper, and J. Graham, "Active shape models—their training and application," *Comput. Vis. Image Underst.*, 1995.

[51] P. Sauer, T. Cootes, and C. Taylor, "Accurate regression procedures for active appearance models," *Proceedings of the British Machine Vision Conference (BMVC)*, 2011.

[52] J. Sung and D. Kim, "Adaptive active appearance model with incremental learning," *Pattern Recogn. Lett.*, 2009.

[53] J. Saragih and R. Goecke, "A Nonlinear Discriminative Approach to AAM Fitting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2007.

[54] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic, "Incremental face alignment in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[55] P. Welinder, P. Perona, and P. Dollar, "Cascaded pose regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.

[56] Z. Feng, J. Kittler, M. Awais, P. Huber, and X. Wu, "Wing loss for robust facial landmark localisation with convolutional neural networks,"

[57] L. Liu, G. Li, Y. Xie, Y. Yu, Q. Wang, and L. Lin, "Facial landmark machines: A backbone-branches architecture with progressive representation learning," *IEEE Transactions on Multimedia*, 2019.

[58] X. Peng, X. Feris, Rogerio S.and Wang, and D. N. Metaxas, "A recurrent encoder-decoder network for sequential face alignment," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.

[59] H. Zhang, Q. Li, and Z. Sun, "Joint voxel and coordinate regression for accurate 3d facial landmark localization," in *International Conference on Pattern Recognition (ICPR)*, 2018.

[60] A. Bulat and G. Tzimiropoulos, "How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks)," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[61] A. Bulat and Y. Tzimiropoulos, "Convolutional aggregation of local evidence for large pose face alignment," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2016.

[62] A. Bulat and G. Tzimiropoulos, "Two-stage convolutional part heatmap regression for the 1st 3d face alignment in the wild (3dfaw) challenge"," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2016.

[63] J. McDonagh and G. Tzimiropoulos, "Joint face detection and alignment with a deformable hough transform model," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2016.

[64] Z. Sánta and Z. Kato, "3d face alignment without correspondences," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2016.

[65] Y. Liu, A. Jourabloo, W. Ren, and X. Liu, "Dense face alignment," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, 2017.

[66] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[67] P. Goyal, P. Dollár, R. B. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch SGD: training imagenet in 1 hour," *arXiv preprint arXiv:1706.02677*, 2017.

[68] C. Bhagavatula, C. Zhu, K. Luu, and M. Savvides, "Faster than real-time facial alignment: A 3d spatial transformer network approach in unconstrained poses," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[69] T.-Y. Yang, Y.-T. Chen, Y.-Y. Lin, and Y.-Y. Chuang, "Fsa-net: Learning fine-grained structure aggregation for head pose estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[70] Z. Liu, Z. Chen, J. Bai, S. Li, and S. Lian, "Facial pose estimation by deep learning from label distributions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2019.

[71] H.-W. Hsu, T.-Y. Wu, S. Wan, W. H. Wong, and C.-Y. Lee, "Quatnet: Quaternion-based head pose estimation with multiregression loss," *IEEE Transactions on Multimedia*, 2019.

[72] H. Zhang, M. Wang, Y. Liu, and Y. Yuan, "Fdn: Feature decoupling network for head pose estimation," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.

[73] R. Valle, J. M. Buenaposada, and L. Baumela, "Multi-task head pose estimation in-the-wild," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[74] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, "Pixel2mesh: Generating 3d mesh models from single rgb images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.