# Deep Face Feature for Face Alignment

Boyi Jiang, Juyong Zhang, Bailin Deng, Yudong Guo and Ligang Liu

*Abstract*—In this paper, we present a deep learning based image feature extraction method designed specifically for face images. To train the feature extraction model, we construct a large scale photo-realistic face image dataset with ground-truth correspondence between multi-view face images, which are synthesized from real photographs via an inverse rendering procedure. The deep face feature (DFF) is trained using correspondence between face images rendered from different views. Using the trained DFF model, we can extract a feature vector for each pixel of a face image, which distinguishes different facial regions and is shown to be more effective than general-purpose feature descriptors for face-related tasks such as matching and alignment. Based on the DFF, we develop a robust face alignment method, which iteratively updates landmarks, pose and 3D shape. Extensive experiments demonstrate that our method can achieve state-of-the-art results for face alignment under highly unconstrained face images.

*Index Terms*—Feature Learning, Face Alignment

## I. INTRODUCTION

Face alignment from images has been an active research topic since 1990s. Face alignment plays an important role in many applications such as 3D face reconstruction [1] and face recognition [2], because it is often used as a pre-processing step. Recently, face alignment has gained significant progress in both theory and practice. Although AAM-based approaches [3], [4], [5], [6] and regression-based approaches [7], [8], [9], [10] work well for face images with small poses, they usually cannot handle profile face images as they do not consider the visibility of landmarks.

In recent years, several methods introduce the 3D Morphable Model (3DMM) [11] for face alignment and achieve better results [12], [1]. Using a 3D face model to compute the visibility and position of 2D landmarks, these methods can handle challenging cases with large pose variation. However, the reconstruction accuracy of such methods is often insufficient. Existing approaches usually use general-purpose features such as SIFT [13] to determine the parameters of a 3DMM by cascaded regression. On the one hand, 3DMM shape and expression parameters are highly non-linear with image texture information, making the mapping difficult to estimate. On the other hand, SIFT type descriptors are designed based on color information of local patches, which works well for general objects but do not make use of specific priors for face images. Therefore, it is possible to achieve better alignment performance by designing a feature descriptor tailored for face images .

Recently, convolution neural networks (CNN) have been successfully applied to many related tasks with state-of-the-art

Boyi Jiang, Juyong Zhang(Corresponding author), Yudong Guo, Ligang Liu are with School of Mathematical Sciences, University of Science and Technology of China. E-mail: jby1993@mail.ustc.edu.cn, juyong@ustc.edu.cn, gyd2011@mail.ustc.edu.cn, lgliu@ustc.edu.cn.

Bailin Deng is with School of Computer Science and Informatics, Cardiff University. E-mail: DengB3@cardiff.ac.uk.

results. In [14], [15], [16], CNN are applied to improve the regression accuracy. In this paper, we propose to use CNN on face images to train a feature extraction model. For each pixel of a face image, we can use this model to extract a high-dimensional feature, which can accurately indicate the same anatomical facial point across different face images under unconstrained conditions, with better performance on many face related tasks than classical features. To train the model, we need a large number of unconstrained face images with registered ground-truth 3DMM faces. However, it is not easy to obtain real face images with ground-truth 3D shapes, especially for profile view face images. To solve this problem, we synthesize a large-scale face image dataset with different poses and expressions together with ground-truth 3D shapes. With the well constructed training set, a novel feature learning method is proposed to extract the feature vector of each pixel such that the feature vector is distinguishable for each face part and smooth over the whole face area. Based on the trained face feature, we design a simple cascaded regression method to perform face alignment in the wild. The experimental results demonstrate that the trained deep face feature has good performance on feature matching, and the face alignment method outperforms existing methods especially on face images with large pose. In summary, the main contributions of this work include:

- We propose a CNN-based face image feature extraction method using a well-constructed training dataset and a novel loss function; the trained feature outperforms general-purpose feature descriptors such as SIFT.
- With the newly designed face feature extractor, we propose a simple yet effective cascaded regression-based approach for face alignment in the wild, and the experimental results demonstrate its better performances over existing methods.

## II. RELATED WORKS

**Classical Face Alignment.** Classical face alignment methods, including Active Shape Model (ASM) [17], [18] and Active Appearance Model (AAM) [5], [4], [3], [6], simulate the image generation process and perform face alignment by minimizing the difference between the model appearance and the input image. These methods can achieve accurate reconstruction results, but require a large number of face models with detailed and accurate point-wise correspondence, as well as high computation cost of parameter fitting. Constrained Local Model (CLM) [19], [20] employs discriminative local texture models to regularize the landmark locations. The CLM algorithm is more robust than the AAM method, which updates the model parameters by minimizing the image reconstruction error. Recently, regression based methods [21], [7] have been proposed to directly estimate landmark locations from the discriminative features around landmarks. Most regression
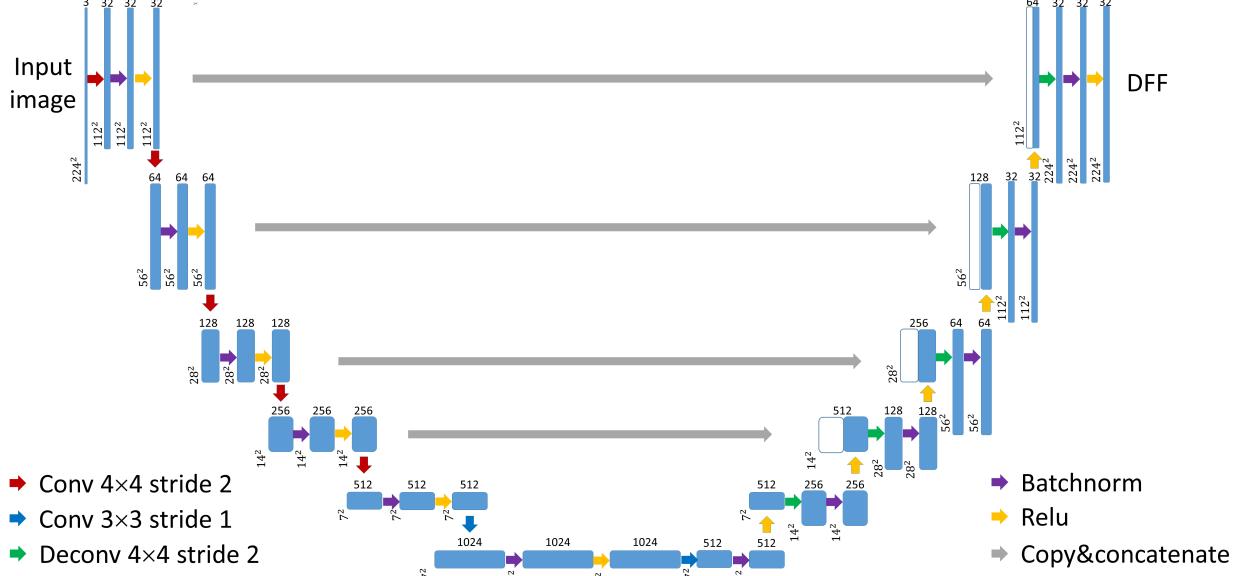
Figure 1: The neural network architecture of DFF extractor. Each blue box represents a multi-channel feature map. The number of channels is on top of the box. The row and column numbers are provided at the lower left edge of the box. White boxes represent copied feature maps.

based algorithms do not consider the visibility of facial landmarks under different view angles. As a result, their performance can degrade substantially for input face images with large poses.

**Face Alignment via Deep Learning.** In recent years, deep learning based methods have been successfylly applied to face alignment and achieved remarkable results. The methods of [22], [23], [24] use multi-stage CNN to regress sparse face image landmark locations. To boost the performance of face alignment, Zhang et al. [25] combine face detection, face alignment, and other tasks into the training of CNN. Alignment of faces with large pose variation is a very challenging problem, because each face might have a different set of visible landmarks. Early works for large-pose face alignment rely on MVS based methods [26], [27], which use different landmark templates for different views and incur high computation costs. Recently, 3D model based techniques [28], [12], [15], [14] have been proposed to address the problem of alignment accuracy for challenging inputs, e.g. those with non-frontal face poses, low image quality, and occlusion, etc. These techniques utilize a 3D morphable model [11], [29] (3DMM) to handle self-occluded landmarks and large-pose landmark detection. Jourabloo and Liu [15] integrate 2D landmark estimation into the 3D face model fitting process, use cascaded CNN to replace simple regressor, and are able to detect 34 landmarks under all view angles. Zhu et al. [14] reduce the CNN regression complexity by estimating a Projected Normalized Coordinate Code map, which can detect 68 anatomical landmarks with visibility judgement. Yu et al. [16] predict dense facial correspondences by training a encoder-decoder network with synthesized face images, and produce robust face image alignment results.

**3D Face Reconstruction.** The 3DMM establishes statistical linear parametric models for both texture and shapes of human faces, with a 3D face represented by a set of coefficients for its shape and texture basis. To recover the face from a 2D image, 3DMM-based methods [11], [30], [29] estimate the shape and texture coefficients by maximizing the similarity between the input 2D face image and the projected 3D face. However, such methods are not robust enough to handle facial landmarks under large pose variation. Multi-view stereo (MVS) [31], [32] is a classical reconstruction method that requires dense correspondence between neighboring images to achieve satisfactory results. When such methods are applied on multi-view face images, the reconstructed face point cloud might contain holes due to insufficient detected matched points.

## III. OUR METHOD

To train the DFF, we first build a large-scale training dataset, which consists of face images and their corresponding ground-truth face shapes and camera parameters. Details on the construction of such training data are provided in Sec. III-A. The DFF is trained using a convolutional neural network (CNN), which will be discussed in Sec. III-B. Finally, we propose a new face alignment algorithm based on the DFF descriptor in Sec. III-C. In the following, we first introduce the 3D face shape representation and the 3D Morphable Model, based on which our algorithm is developed.

**3D Face Shape.** We represent a 3D face shape using a triangle mesh with $n$ vertices and fixed connectivity. Therefore, each face shape is determined by its vertex positions, represented as a $3 \times n$ matrix

$$\mathbf{S} = \begin{pmatrix} x_1 & x_2 & \cdots & x_n \\ y_1 & y_2 & \cdots & y_n \\ z_1 & z_2 & \cdots & z_n \end{pmatrix}, \tag{1}$$

where each column vector corresponds to the 3D coordinates of a vertex. We denote a face image as $\mathbf{I}$, and assume the mapping from $\mathbf{S}$ to $\mathbf{I}$ to be a weak perspective projection with camera

Figure 2: Examples of augmented training images. From the original face image (left), we construct a set of new images with different view directions and expressions.

parameters $\mathbf{w} = (s, \alpha, \beta, \gamma, t_x, t_y)$, where $(t_x, t_y)$ represent the translation on the image plane, $s$ the scaling factor, $\alpha$ the pitch angle, $\beta$ the yaw angle, and $\gamma$ the roll angle. The three angles $(\alpha, \beta, \gamma)$ determine a rotation matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$. Then the 2D projection of a vertex $\mathbf{q} \in \mathbb{R}^3$ from the 3D face model can be written as:

$$P(\mathbf{q}) = s \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \mathbf{R}\mathbf{q} + \mathbf{t}, \qquad (2)$$

where $\mathbf{t} = \begin{pmatrix} t_x & t_y \end{pmatrix}^T$.

**Parametric Face Model.** The excessive degrees of freedom from vertex positions in Equation (1) can result in unnatural face shapes. From a perceptual point of view, realistic face shapes lie in a space of lower dimension [33]. Therefore, we follow the approach of 3D Morphable Model (3DMM) [11] and represent a face shape and its albedo using a lower-dimensional linear model. We also introduce delta blendshapes to represent the deformation from a neutral face to its different expressions [34]. The resulting parametric face model is

$$\begin{aligned} \widehat{\mathbf{S}} &= \overline{\mathbf{S}} + \mathbf{A}_{\text{id}}\mathbf{p}_{\text{id}} + \mathbf{A}_{\text{exp}}\mathbf{p}_{\text{exp}}, \\ \widehat{\mathbf{T}} &= \overline{\mathbf{T}} + \mathbf{B}_{\text{alb}}\mathbf{p}_{\text{alb}}, \end{aligned} \qquad (3)$$

where $\widehat{\mathbf{S}} \in \mathbb{R}^{3n}$ stacks the vertex coordinates of $\mathbf{S}$, and $\widehat{\mathbf{T}} \in \mathbb{R}^{3n}$ represent the albedo values for the vertices. Here $\overline{\mathbf{S}} = \overline{\mathbf{S}}_{\text{id}} + \overline{\mathbf{S}}_{\text{exp}} \in \mathbf{R}^{3n}$ is the mean face shape, where $\overline{\mathbf{S}}_{\text{id}}$ and $\overline{\mathbf{S}}_{\text{exp}}$ are the mean identity and mean expression; $\overline{\mathbf{T}} \in \mathbf{R}^{3n}$ denotes the mean albedo values; $\mathbf{A}_{\text{id}}, \mathbf{B}_{\text{alb}} \in \mathbf{R}^{3n \times m_1}$ and $\mathbf{A}_{\text{exp}} \in \mathbb{R}^{3n \times m_2}$ are the bases for identity, albedo and expression, respectively, computed using PCA; $\mathbf{p}_{\text{id}}, \mathbf{p}_{\text{exp}}, \mathbf{p}_{\text{al}}$ are their linear combination coefficients. We denote the collection of shape parameters by $\mathbf{p} = (\mathbf{p}_{\text{id}}^T, \mathbf{p}_{\text{exp}}^T)^T$. We choose $m_1 = 80$ and $m_2 = 79$ for our implementation. The mean and basis identities are constructed from the Basel Face Model [35], while the mean and basis expressions are constructed using the FaceWarehouse dataset [36].

*A. Training Data Construction*

To train the feature extraction model DFF, we need a large set of training data $\{(\mathbf{S}_i, \mathbf{I}_i, \mathbf{w}_i) \mid i = 1, \cdots, N\}$, where $\mathbf{I}_i$ is the $i$-th face image, and $\mathbf{S}_i, \mathbf{w}_i$ are the corresponding ground-truth face shape and camera parameters. From $\mathbf{S}_i$ and $\mathbf{w}_i$, we can easily compute the location $\mathbf{U}_i$ and visibility $\mathbf{V}_i$ of ground-truth landmarks, by projecting the pre-selected landmark vertices from the face shape $\mathbf{S}_i$ onto the image plane and checking their visibility in 3D.

We first select 4308 face images from the 300W dataset [37] and the Multi-PIE dataset [38], and then follow the approach

of [39] to fit the parametric face model to each image and produce the shape, albedo, and camera parameters. Using these data, new face images are then generated by rendering the face shape with various expressions and camera parameters, resulting in a dataset that includes 80000 face images with synthesized ground-truth face shapes and camera parameters. Each face image has a resolution of $224 \times 224$. Fig. 2 shows an example of the generated images.
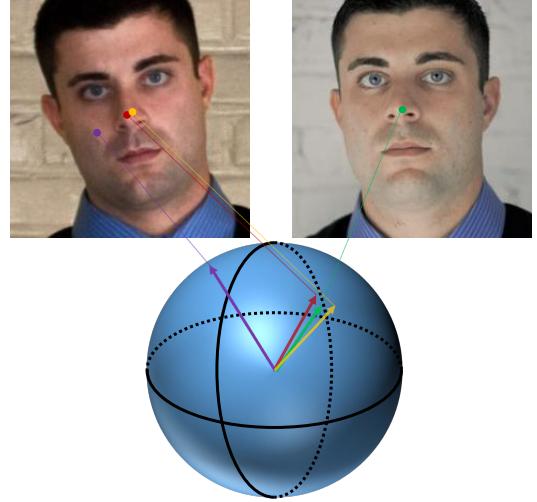


Figure 3: An illustration of the requirements for the DFF. If two pixels correspond to nearby points on the 3D face surface, their normalized DFF's should be close on the hypersphere; otherwise, they should be sufficiently far away from each other.

*B. Deep Face Feature Training*

Many existing face alignment methods [9], [40] rely on local features such as local binary feature (LBF) and SIFT. Although these methods work well in many scenarios, they might fail for input face images with large poses. One potential cause of the problem is that local features cannot utilize global information from the whole image and fail to recognize global structures such as specific face regions and self-occlusion. In this work, we propose a deep learning based end-to-end method, to extract for each face image pixel a feature vector that takes global information into account. Our method uses a neural network to map each pixel to a high-dimensional point, which is then normalized to have unit length. To effectively indicate and distinguish facial features, the normalized DFF descriptor should preserve the metric structure on the 3D face surface. In particular, for two pixels corresponding to the same anatomical region, their normalized DFF descriptors should be close to each other even if they are from different images with different poses, scales, and lighting conditions. On the other hand, for pixels corresponding to different facial parts, their normalized DFF descriptors should be sufficiently far away from each other even if their surrounding image regions have similar appearances. One such example is shown in Fig. 3.

To satisfy these criteria, we follow the approach of [41] and train the DFF extractor to solve a series of classification problems. The approach is based on the following observation:

Figure 4: Examples of random segmentations of the 3D face surface into 500 patches.



Figure 5: For each image corresponding to a segmented 3D face model, the visible patches are projected back onto the image plane to label the pixels they cover.

if we randomly segment the face surface into uniform patches, then nearby points on the face surface are likely to lie on the same patch. Accordingly, if we perform classification of face image pixels into the segmentation patches according to their DFF descriptors, then pixels corresponding to nearby 3D face points should be classified into the same patch with high probability. In other words, the DFF extractor should lead to a small value of the classification loss function corresponding to the segmentation. To avoid bias towards a specific segmentation, we generate a large set of random segmentations for each face, and use the sum of their classification loss functions as the overall loss function for training the DFF extractor.

In detail, we randomly generate 100 uniform segmentations of the mean face model, where each segmentation consists of 500 patches, and each set is a union of mesh faces (see Fig. 4 for examples). The segmentation is performed by computing centroidal Voronoi tessellation [42] on the surface, using random sample points as initial generators. Then for each ground-truth 3D face $\mathbf{S}_i$ in the training data set, we derive 100 segmentations by using the same set of faces for each patch as the mean face segmentation. For each segmentation, we project the visible patches to each training face image corresponding to $\mathbf{S}_i$, with the patch visibility determined from the camera parameters (see Fig. 5 for an example). The image pixels are then labeled according to the projected visible patches. We then use all the images with labels to train a CNN per-pixel classifier. The CNN consists of two parts. The first part is a DFF extractor that takes a $224 \times 224$ color image as input, and produces a 32-dimensional feature vector for each pixel. The architecture of our DFF extractor is similar to the u-net in [43], applying convolution and de-convolution layer symmetrically and concatenating with shallow feature map (see Fig. 1). The second part consists of 100 independent classification loss layers, one for each segmentation. Each layer takes the generated DFF's as input, and evaluates a classification loss function according to the segmentation. The sum of all loss functions is the final loss function for training the network. For each segmentation, we utilize the angular softmax [44] with $m = 1$ as the classification loss, which assumes a feature vector for each class such that input vectors are classified based on their angle to the class feature vectors. For each segmentation, we denote by $\mathbf{h}_j \in \mathbb{R}^{32}$ the unit feature vector for patch $j$ ($j = 1, \ldots, 500$), which is part of the parameters for the classification loss layer. Then its classification loss for

an image is written as

$$l_k = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} -\log\left(\frac{\exp(\mathbf{h}_{\tau(p)} \cdot \mathbf{f}_p)}{\sum_{j=1}^{500} \exp(\mathbf{h}_j \cdot \mathbf{f}_p)}\right),$$

where $\mathcal{P}$ is the set of pixels corresponding to the visible patches, $\tau(p)$ is the index of the patch corresponding to pixel $p$, and $\mathbf{f}_p$ is the normalized DFF descriptor for $p$. The training optimizes the parameters of the DFF extractor as well as the classification loss layers, to reduce the sum of all classification loss.

*C. Large Pose Face Alignment*

Classical cascaded regression algorithms such as the Supervised Descent Method (SDM) [9] use SIFT descriptors around the landmarks to regress their locations. These methods achieve state-of-the-art performance for small- and medium-pose face alignment. However, their results are not satisfactory for large poses or more complicated scene conditions. This is because for general descriptors such as SIFT, their capability of identifying key points can decrease drastically under large pose variation. Moreover, without visibility information, SIFT descriptors for landmarks that are occluded in some images may interfere with the regression.

To address these issues, we propose an iterative approach that utilizes the DFF descriptor instead of SIFT, and uses the parametric face model to estimate face pose and landmark visibility. Given an input face image, we find a 3D parametric face model as well as the camera parameters, such that the projected 3D landmarks according to the camera parameters align with the face image. Starting from the initial parametric face model and the initial camera parameters, we first project the visible 3D landmark vertices onto the image plane to obtain 2D landmark locations, and evaluate the DFF descriptors of the face image at these landmark locations. From the DFF descriptors, we determine the target new locations of 2D landmarks and update the camera parameters, using a learned mapping from DFF descriptors to landmark and camera parameter updates that improve their accuracy. Afterwards, the parametric 3D face model is updated accordingly, such that its
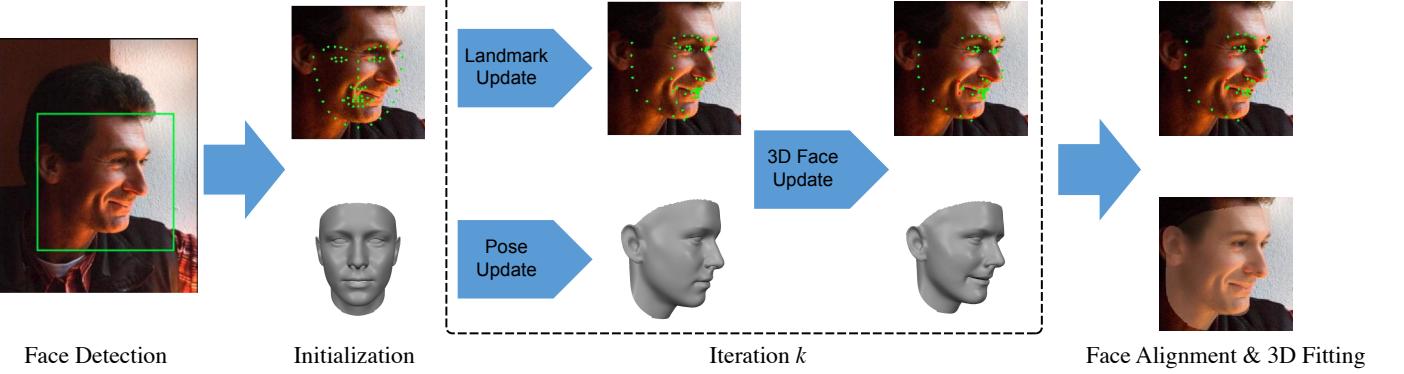
Figure 6: An overview of our face alignment algorithm pipeline. The input is a face image with a detected face region. Starting from the initial 3D face model and its projected 2D landmark locations (shown in green), we evaluate the DFF descriptors at the projected landmarks, and use them to compute the target landmark locations and update the camera parameters according to a generic descent direction learned from the training data. Afterwards, we update the 3D face model to align with the target landmark locations, and recompute its projected 2D landmarks and their visibility (with invisible landmarks shown in red). This process is iterated until convergence.

projected landmark locations using the new camera parameters are as close as possible to the target locations. This process is iterated until convergence. Figure 6 illustrates the pipeline of our algorithm. In the following, we explain each step in detail.

*1) Initialization:* Given an input face image, we first run a face detector to locate the face region. We initialize the parametric face model $\mathbf{S}^{(0)}$ and the camera parameters $\mathbf{w}^{(0)}$ to the mean face shape and mean camera parameters, respectively. The landmark vertices on $\mathbf{S}^{(0)}$ are projected onto the image plane according to Equation (2) using camera parameters $\mathbf{w}^{(0)}$, to determine the initial 2D landmark locations as well as their visibility. In this work, the alignment between the 3D face model and the 2D face image is done via 68 landmark vertices, chosen on the 3D face mesh to match the 68 landmark points used in the Multi-PIE data set [38] (see Fig. 7 left). The initial 2D locations of these landmarks are denoted by $\mathbf{X}^{(0)} \in \mathbb{R}^{136}$. In addition, we use a denser set of landmark vertices to evaluate DFF descriptors, in order to capture more global information of the face shape and produce more robust results. Specifically, we
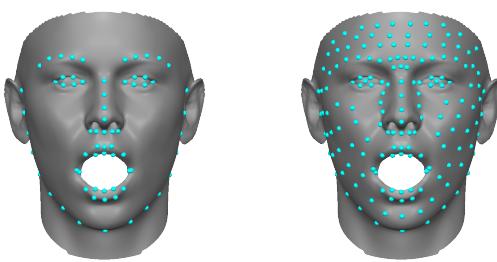


Figure 7: We perform face alignment using 68 landmarks (shown on the left), which are chosen on the 3D face mesh according to the 68 feature points of the Multi-PIE dataset [38]. In our face alignment algorithm, the update of landmarks and camera parameters is driven by DFF descriptors corresponding to 160 landmarks on the 3D face model (shown on the right), which is a superset of the 68 landmarks.

densely sample the face mesh to obtain 160 landmark vertices, which includes the previous 68 landmarks (see Fig. 7 right). The initial 2D locations of these landmarks are denoted by $\mathbf{U}^{(0)} \in \mathbb{R}^{320}$, and their visibility is indicated using a binary vector $\mathbf{V}^{(0)} \in \mathbb{R}^{160}$.

*2) Updating Landmarks and Camera Parameters:* In the $k$-th iteration, we first evaluate the DFF descriptors of the dense landmark set using their current locations $\mathbf{U}^{(k)}$ and visibility $\mathbf{V}^{(k)}$. The DFF descriptors are concatenated into a vector $\mathbf{F}^{(k)} \in \mathbb{R}^{5120}$. For invisible landmarks, their corresponding components $\mathbf{F}^{(k)}$ are set to zero. Using $\mathbf{F}^{(k)}$, we compute new camera parameters $\mathbf{w}^{(k+1)}$ as well as target locations for new 2D landmarks $\widehat{\mathbf{X}}^{(k+1)}$, both of which should improve the accuracy of alignment. We adapt the approach of [9] to determine $\mathbf{w}^{(k+1)}$ and $\widehat{\mathbf{X}}^{(k+1)}$, by computing the displacements $\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}$ and $\widehat{\mathbf{X}}^{(k+1)} - \mathbf{X}^{(k)}$ as a linear function of the DFF descriptors $\mathbf{F}^{(k)}$:

$$\widehat{\mathbf{X}}^{(k+1)} = \mathbf{X}^{(k)} + \mathbf{R}_{\mathbf{X}}^{(k)}\mathbf{F}^{(k)} + \mathbf{b}_{\mathbf{X}}^{(k)}, \quad (4)$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{R}_{\mathbf{w}}^{(k)}\mathbf{F}^{(k)} + \mathbf{b}_{\mathbf{w}}^{(k)}. \quad (5)$$

Here matrix $\mathbf{R}_{\mathbf{X}}^{(k)}$ and vector $\mathbf{b}_{\mathbf{X}}^{(k)}$ represent a *generic descent direction* that improves the accuracy of landmark locations according to the current DFF descriptors [9]. They are learned from the training data set constructed in Sec. III-A, by solving a sequence of regression problems. Similarly, matrix $\mathbf{R}_{\mathbf{w}}^{(k)}$ and vector $\mathbf{b}_{\mathbf{w}}^{(k)}$ are learned from the training data set. The details of learning the generic descent directions are explained later in Sec. III-C4.

*3) Updating 3D Face Model:* Since the target 2D landmark positions $\widehat{\mathbf{X}}^{(k+1)}$ are determined by simply applying a generic descent step, they do not necessarily correspond to a specific parametric 3D face model. Therefore, to ensure the validity of landmark locations, we compute an updated parametric face model $\mathbf{S}^{(k+1)}$ according to $\widehat{\mathbf{X}}^{(k+1)}$, and determine new 2D landmark locations by projecting the updated face model according to the camera parameters $\mathbf{w}^{(k+1)}$. This results in an optimization problem for the shape parameters $\mathbf{p} =$

$(\mathbf{p}_{\mathrm{id}}^T, \mathbf{p}_{\mathrm{exp}}^T)^T$ that determine the face model $\mathbf{S}$ according to Equation (3):

$$\mathbf{p}^{(k+1)} = \underset{\mathbf{p}}{\arg\min}\ \omega_{\mathrm{lan}} E_{\mathrm{lan}}^{(k)}(\mathbf{p}) + \omega_{\mathrm{reg}} E_{\mathrm{reg}}^{(k)}(\mathbf{p}). \qquad (6)$$

Here $E_{\mathrm{lan}}^{(k)}$ is a landmark fitting term, $E_{\mathrm{reg}}^{(k)}$ is a regularization term, and $\omega_{\mathrm{lan}}, \omega_{\mathrm{reg}}$ are their weights. The terms are defined as:

$$E_{\mathrm{lan}}^{(k)}(\mathbf{p}) = \|\widehat{\mathbf{X}}^{(k+1)} - \mathbf{Y}(\mathbf{w}^{(k+1)}, \mathbf{p})\|^2, \qquad (7)$$

$$E_{\mathrm{reg}}^{(k)}(\mathbf{p}) = \mathbf{p}_{\mathrm{id}}^T\ \mathbf{D}_{\mathrm{id}}^{-1}\ \mathbf{p}_{\mathrm{id}} + \mathbf{p}_{\mathrm{exp}}^T\ \mathbf{D}_{\mathrm{exp}}^{-1}\ \mathbf{p}_{\mathrm{exp}}. \qquad (8)$$

Here $\mathbf{Y}(\mathbf{w}^{(k+1)}, \mathbf{p})$ is a vector storing the projected 2D coordinates of the 68 landmarks from the face model with shape parameters $\mathbf{p}$, according to camera parameters $\mathbf{w}^{(k+1)}$. $\mathbf{D}_{\mathrm{id}}$ and $\mathbf{D}_{\mathrm{exp}}$ are diagonal matrices that store the eigenvalues of the identity and expression covariance matrices corresponding to bases $\mathbf{A}_{\mathrm{id}}$ and $\mathbf{A}_{\mathrm{exp}}$ in Equation (3). This is a linear least-squares problem and can be solved efficiently.

After solving for the shape parameters $\mathbf{p}^{(k+1)}$, we determine the new 3D face shape $\mathbf{S}^{(k+1)}$, and project its landmark vertices to obtain the new 2D landmark coordinates $\mathbf{U}^{(k+1)}, \mathbf{X}^{(k+1)}$, as well as the visibility $\mathbf{V}^{(k+1)}$ for $\mathbf{U}^{(k+1)}$. $\mathbf{V}^{(k+1)}$ can be easily determined on the GPU by rendering the triangles of the 3D face mesh and the landmark vertices to the OpenGL depth buffer, with the OpenGL projection matrix determined from the corresponding camera parameters.

We iteratively apply landmark update, camera parameters update and 3D face update until convergence. In our experiments, three iterations are sufficient for good results.

*4) Learning the Generic Descent Directions:* In Sec. III-C2, the computation of target landmark positions $\widehat{\mathbf{X}}^{(k+1)}$ and new camera parameters $\mathbf{w}^{(k+1)}$ require the generic descent directions — represented using matrices $\mathbf{R}_{\mathbf{X}}^{(k)}, \mathbf{R}_{\mathbf{w}}^{(k)}$ and vectors $\mathbf{b}_{\mathbf{X}}^{(k)}, \mathbf{b}_{\mathbf{w}}^{(k)}$ — that map the current DFF descriptors to the displacements from landmark positions $\mathbf{X}^{(k)}$ and camera parameters $\mathbf{w}^{(k)}$, respectively. Following [9], we learn the generic descent directions in each iteration using the training images $\mathbf{I}_i$ ($i = 1, \ldots, N$) together with their ground-truth landmark locations $\mathbf{X}_i^*$ and camera parameters $\mathbf{w}_i^*$. Using a generic face model $\overline{\mathbf{S}}$ for all the training images, we iteratively perform landmark update, camera parameters update and 3D face update similar to Secs. III-C1 to III-C3, and determine the generic descent directions via linear regression.

In detail, we initialize the $\overline{\mathbf{S}}^{(0)}$ to the mean face shape, and initialize the camera parameters $\mathbf{w}_i^{(0)}$, landmark locations $\mathbf{X}_i^{(0)}, \mathbf{U}_i^{(0)}$ and visibility $\mathbf{V}_i^{(0)}$ for image $\mathbf{I}_i$ as described in Sec. III-C1. In the $k$-th iteration, we evaluate the DFF descriptors $\mathbf{F}_i^{(k)}$ of image $\mathbf{I}_i$ according to its landmark locations $\mathbf{U}_i^{(k)}$ and visibility $\mathbf{V}_i^{(k)}$, as described in Sec. (III-C2). The generic descent direction $(\mathbf{R}_{\mathbf{X}}^{(k)}, \mathbf{b}_{\mathbf{X}}^{(k)})$ for landmark locations in iteration $k$ is then determined such that the resulting updated landmarks are as close as possible to the ground-truth landmarks $\mathbf{X}_i^*$ for all the training images. We thus compute

$(\mathbf{R}_{\mathbf{X}}^{(k)}, \mathbf{b}_{\mathbf{X}}^{(k)})$ by solving a regression problem

$$(\mathbf{R}_{\mathbf{X}}^{(k)}, \mathbf{b}_{\mathbf{X}}^{(k)}) = \min_{\mathbf{R}, \mathbf{b}} \sum_{i=1}^{N} \|\mathbf{X}_i^* - \mathbf{X}_i^{(k)} - \mathbf{R}\mathbf{F}_i^{(k)} - \mathbf{b}\|_2^2 \\ + \lambda_1 \left(\|\mathbf{R}\|_F^2 + \|\mathbf{b}\|_2^2\right), \qquad (9)$$

where the second term is a regularization, and $\lambda_1$ is a positive weight. Similarly, the generic descent direction $(\mathbf{R}_{\mathbf{w}}^{(k)}, \mathbf{b}_{\mathbf{w}}^{(k)})$ is determined via a regression problem

$$(\mathbf{R}_{\mathbf{w}}^{(k)}, \mathbf{b}_{\mathbf{w}}^{(k)}) = \min_{\mathbf{R}, \mathbf{b}} \sum_{i=1}^{N} \|\mathbf{w}_i^* - \mathbf{w}_i^{(k)} - \mathbf{R}\mathbf{F}_i^{(k)} - \mathbf{b}\|_2^2 \\ + \lambda_2 \left(\|\mathbf{R}\|_F^2 + \|\mathbf{b}\|_2^2\right). \qquad (10)$$

Both regressions are linear least-squares problems and can be solved efficiently. Afterwards, we compute the target landmark locations $\widehat{\mathbf{X}}_i^{(k+1)}$ and updated camera parameters $\mathbf{w}_i^{(k+1)}$ for each training image according to Equations (4) and (5). Then the generic face model $\overline{\mathbf{S}}$ is updated by optimizing its shape parameters $\overline{\mathbf{p}} = (\overline{\mathbf{p}}_{\mathrm{id}}^T, \overline{\mathbf{p}}_{\mathrm{exp}}^T)^T$ to align with the target landmark locations in all training images, similar to Equation (6). To be precise, we solve for

$$\overline{\mathbf{p}}^{(k+1)} = \underset{\overline{\mathbf{p}}}{\arg\min} \frac{\omega_{\mathrm{lan}}}{N} \sum_{i=1}^{N} \|\widehat{\mathbf{X}}_i^{(k+1)} - \mathbf{Y}(\mathbf{w}_i^{(k+1)}, \overline{\mathbf{p}})\|^2 \\ + \omega_{\mathrm{reg}} \left(\overline{\mathbf{p}}_{\mathrm{id}}^T\ \mathbf{D}_{\mathrm{id}}^{-1}\ \overline{\mathbf{p}}_{\mathrm{id}} + \overline{\mathbf{p}}_{\mathrm{exp}}^T\ \mathbf{D}_{\mathrm{exp}}^{-1}\ \overline{\mathbf{p}}_{\mathrm{exp}}\right). \qquad (11)$$

Using $\overline{\mathbf{p}}^{(k+1)}$, we derive the new generic face model $\overline{\mathbf{S}}^{(k+1)}$, and compute the new landmark locations $\mathbf{X}_i^{(k+1)}, \mathbf{U}_i^{(k+1)}$ and visibility $\mathbf{V}_i^{(k+1)}$ for each training image. We repeat the above procedures until the generic descent directions are learned for all the required iterations.

## IV. EXPERIMENTS

In this section, we evaluate the effectiveness of our approach by applying it to feature point matching and face alignment, and comparing the results with other existing approaches.

### A. Feature Matching

We first evaluate the effectiveness of DFF in capturing global structural information and identifying facial features. Since SIFT is widely used for feature matching and landmark regression, we compare the performance of DFF and SIFT in matching features of multi-view face images. For SIFT matching, we follow the approach of [13]. Specifically, given the source image $\overline{\mathbf{I}}_1$ and the target image $\overline{\mathbf{I}}_2$, we first identify for each image a set of feature points as describe in [13]. We denote the two feature point sets by $S_1$ and $S_2$, respectively. For each point $p_1$ in the set $S_1$, we find a point $p_2$ from the set $S_2$ whose SIFT descriptor is the closest to that of $p_1$. Following [13], we consider the pair $(p_1, p_2)$ is as valid matching only if the ratio $d_{1,2}/d'_{1,2}$ is smaller than a certain threshold, where $d_{1,2}, d'_{1,2}$ are the SIFT descriptor angles between $p_1, p_2$ and between $p_1, p'_2$ respectively, with $p'_2$ a point from $S_2$ whose SIFT descriptor is the second closest to that of $p_1$. In our experiments, the ratio threshold is fine-tuned to $1/1.3$ to achieve the best results.
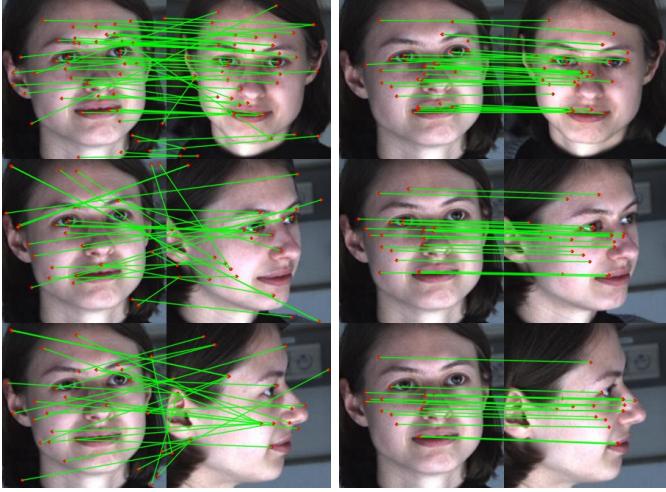
Figure 8: Comparison of feature matching results using SIFT (left) and our DFF descriptor (right), on three image pairs that share the same target image and use different target images with small, medium, and large view angle differences.

For DFF matching, we first evaluate the DFF descriptor for each feature point from $S_1$ that lies on the face region, as well as DFF descriptors for all the face pixels in image $\bar{\mathbf{I}}_2$. Then for each face feature point $q_1$ from $S_1$, we find a face point $q_2$ from $\bar{\mathbf{I}}_2$ whose DFF descriptor vector has the smallest angle from the DFF descriptor of $q_1$. For accurate matching, the pair is considered valid only if the angle between their descriptor vectors is less than a certain threshold. We set the threshold to $30°$ in our experiments, to trade off between accuracy and the number of valid pairs.

Fig. 8 compares the results using the two approaches, on three image pairs with the same source image and different target images with small, medium, and large view angle difference, respectively. It can be seen that DFF leads to more stable and accurate than SIFT, especially for the image pair with large view angle difference. Moreover, since the DFF descriptors encode global structural information, the DFF matching results are consistent across different image pairs, with the same source feature point mapped to consistent points in different target images. This is not the case for SIFT, as it only considers local features. This is apparent on the image pair with large view angle difference, where DFF matching excludes source feature points that corresponding to invisible regions in the target image, while SIFT matches them to the other side of the face due to local structural similarity.

We also apply DFF for dense matching between two face images with different poses, using images from the AFLW2000-3D dataset [14]. We evaluate the DFF descriptor for each face pixel in the source image as well as the target image. Then each source face pixel is matched to a face pixel in the target image whose DFF descriptor has the smallest angle from the source DFF descriptor. The matching is considered as valid if the angle is smaller than $12°$. Here we use a smaller threshold for valid matching, in order to reduce ambiguity due to the large number of source and target pixels. Fig. 9 shows some dense matching results, using color coding to indicate corresponding pixels.
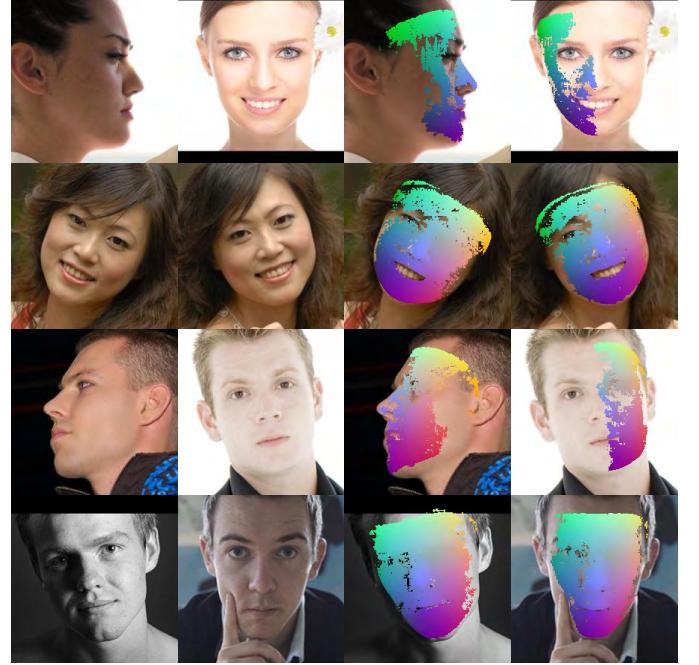


Figure 9: Dense correspondence between multi-view face images based on the DFF descriptor. Each row shows the source image and the target image, and visualizes the correspondence between their face pixels via color coding.

It shows excellent performance of DFF for dense matching between face images with very different poses.

### B. Face Alignment

In Figure 10, we evaluate the robustness of our face alignment algorithm, by applying it to the face images with large area of self shadow and with large face pose. For comparison, we also apply a method similar to the one in Sec. III-C, using SIFT instead of DFF as the feature descriptors. Thanks to global information encoded in DFF, it produces more robust and accurate results than SIFT on these challenging examples.

To evaluate the performance of our approach for large pose face alignment, we test it using face images in the wild from the AFLW dataset[1] and the AFLW2000-3D dataset[2]. The AFLW dataset contains face images with 21 visible ground truth landmarks, while the AFLW2000-3D dataset consists of fitted 3D faces for the first 2000 AFLW samples and can be used for 3D face alignment evaluation. In our experiments, the alignment accuracy is measured by the Normalized Mean Error (NME), which is the average of landmark error normalized by the bounding box size [12]. For AFLW, we compute NME using only the visible landmarks and the bounding boxes provided in the dataset. For AFLW2000-3D, the bounding box that encloses all the 68 ground truth landmarks are used to compute NME, similar to [14]. Since the results reported in [14] are obtained using a model trained with the 300W-LP dataset, we also learn

---

[1]https://www.tugraz.at/institute/icg/research/team-bischof/lrs/downloads/aflw/

[2]http://www.cbsr.ia.ac.cn/users/xiangyuzhu/projects/3DDFA/main.htm

Figure 10: Face alignment results from the method described in Sec. III-C, using SIFT (top row) and DFF (bottom row) as the feature descriptor, respectively. Due to the global information encoded in DFF, it leads to more robust and accurate results than SIFT.

Table II: The NME (%) on AFLW2000-3D Dataset (68 Points)

| Method | [0,30] | [30,60] | [60,90] | Mean | Std |
|---|---|---|---|---|---|
| RCPR [21] (300W-LP) | 4.26 | 5.96 | 13.18 | 7.80 | 4.74 |
| ESR [7] (300W-LP) | 4.60 | 6.7 | 12.67 | 7.99 | 4.19 |
| SDM [9] (300W-LP) | 3.67 | 4.94 | 9.76 | 6.12 | 3.21 |
| 3DDFA [14] | 3.78 | 4.54 | 7.93 | 5.42 | 2.21 |
| 3DDFA+SDM | 3.43 | **4.24** | 7.17 | 4.94 | 1.97 |
| Yu et al. [16] | 3.62 | 6.06 | 9.56 | 6.41 | 2.99 |
| Ours (SIFT) | 5.35 | 6.75 | 8.23 | 6.78 | **1.44** |
| Ours (DFF) | **3.20** | 4.68 | **6.28** | **4.72** | 1.54 |

Table III: The NME (%) on AFLW2000 for Visible Inner Landmarks

| Method | [0,30] | [30,60] | [60,90] | All Images |
|---|---|---|---|---|
| Zhu et al. [14] | 4.30 | 4.41 | 6.68 | 4.60 |
| Yu et al. [16] | 3.14 | 3.84 | 5.53 | 3.58 |
| Ours (DFF) | **2.56** | **3.80** | **4.80** | **3.14** |

the generic descent directions using the 300W-LP dataset for consistency.

All 24384 face images from the AFLW dataset, with yaw angles ranging from $-90°$ to $90°$, are used for testing. To show the robustness of our method for large yaw angles, we divide the images into 3 subsets according to their absolute yaw angles: $[0°, 30°]$, $[30°, 60°]$ and $[60°, 90°]$, with 14032, 5949 and 4403 images respectively. For the AFLW2000-3D dataset, we follow the same experimental setting as [14]. Tables I and II compare the results using our approach and other existing alignment methods on the two datasets, with the best results highlighted in bold font. The results of other methods are gathered from [14] and [16]. Our method significantly outperforms RCPR [21], ESR [7], SDM [9] and Yu et al. [16], especially for samples within $[60°, 90°]$ yaw angles. Our method achieves similar performance with 3DDFA+SDM (which applies SDM to refine the results from 3DDFA [14]) for $[0, 60°]$ yaw angles, and better results for $[60°, 90°]$ yaw angles.

Yu et al. [16] observed that a decent-quality fitting can have a high NME due to the subjective nature of the contour and invisible landmarks in some results of AFLW2000-3D. To better understand the performance of their method, they exclude the contour and invisible landmarks and evaluate the NME using only the inner and visible landmarks. Tab. III shows the comparison result with the same setting, where our method achieves the best results.

Fig. 11 shows some results of our 2D and 3D facial alignment

on challenging and large pose images from AFLW2000. We also show the results using the large pose face alignment methods from [14], [15], and the widely used state-of-the-art face trackers of Kazemi et al. [45]. It can be observed that our method is more robust to heavy occlusions, large variations in illumination, translation, and rotation. More alignment results by our method on AFLW dataset are given in Fig. 12.

We also test the performance of DFF on medium pose face alignment problems. The experiments are conducted on the 300W dataset, and we use the training data of the LFPW and HELEN datasets, and the whole AFW dataset, to train the alignment model. The testing is performed on three parts: the test samples from LFPW and HELEN as the common subset, the 135-image IBUG as the challenging subset, and the union of them as the full set (689 images in total), following the setting of [14]. Since the 300W dataset does not provide ground-truth camera parameters, the large-pose face alignment method described in Sec. III-C can not be used. Instead, we assume all landmarks are visible, and apply the SDM method [9] to regress the landmark positions in a cascaded way, using DFF instead of SIFT as the feature descriptor. Each training image is augmented with 5 random face box, resulting in training images 15740 in total. The alignment accuracy is evaluated by the landmark NME, using the inter-pupil distance for normalization. Tab. IV compares our results with the results reported in [14]. It can be observed that our method greatly

Table I: The NME (%) on AFLW Dataset (21 Points)

| Method | [0,30] | [30,60] | [60,90] | Mean | Std |
|---|---|---|---|---|---|
| RCPR [21] (300W-LP) | 5.43 | 6.58 | 11.53 | 7.85 | 3.24 |
| ESR [7] (300W-LP) | 5.66 | 7.12 | 11.94 | 8.24 | 3.29 |
| SDM [9] (300W-LP) | 4.75 | 5.55 | 9.34 | 6.55 | 2.45 |
| 3DDFA [14] | 5.00 | 5.06 | 6.74 | 5.60 | 0.99 |
| 3DDFA+SDM | 4.75 | **4.83** | 6.38 | 5.32 | **0.92** |
| Yu et al. [16] | 5.94 | 6.48 | 7.96 | 6.79 | 1.05 |
| Ours (SIFT) | 5.65 | 6.23 | 9.24 | 7.04 | 1.93 |
| Ours (DFF) | **3.68** | 5.03 | **5.78** | **4.83** | 1.06 |

Table IV: The NME(%) of Face Alignment Results on 300W, With the First and the Second Best Results Highlighted

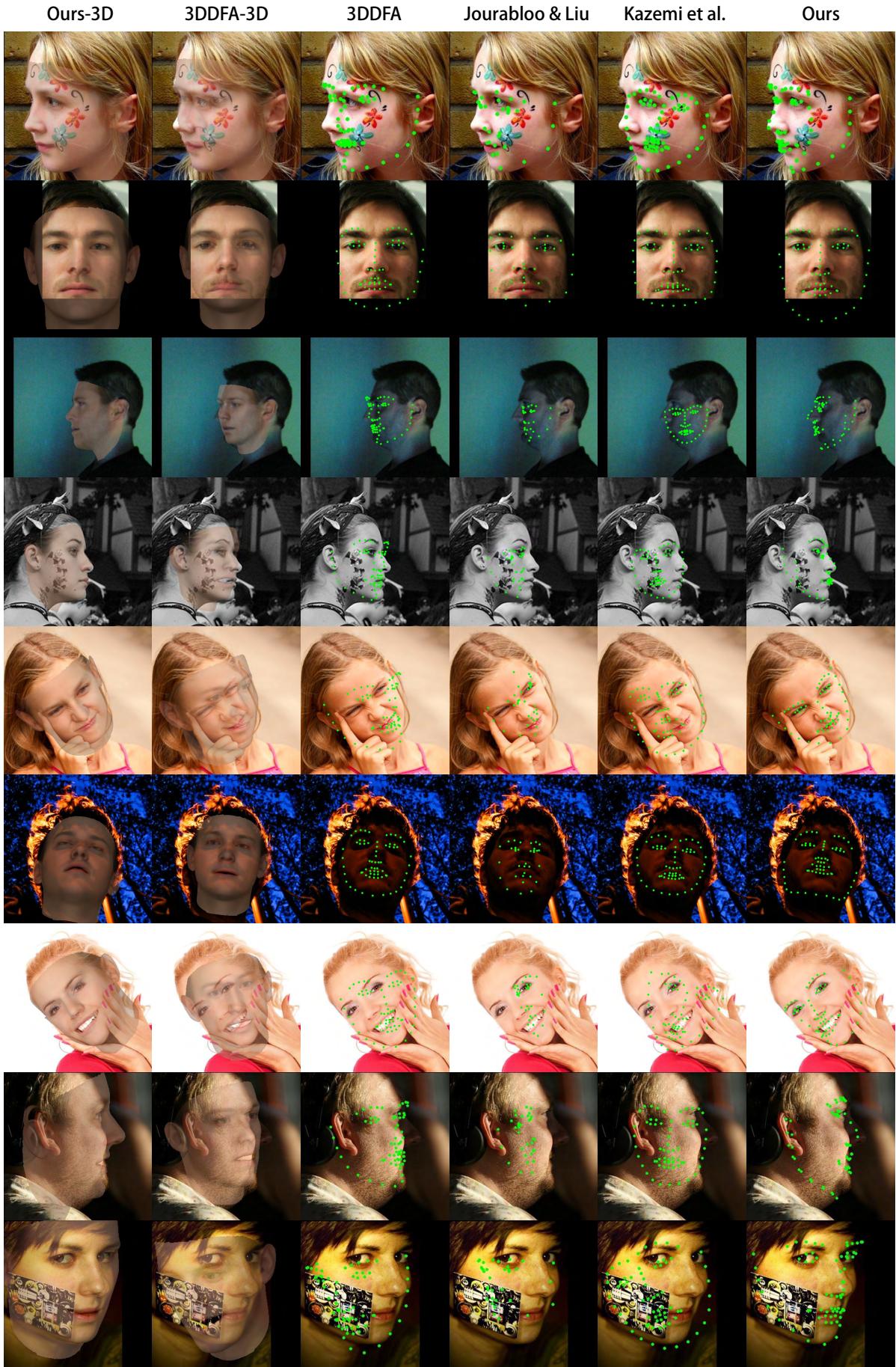| Method | Common | Challenging | Full |
|---|---|---|---|
| TSPM [27] | 8.22 | 18.33 | 10.20 |
| ESR [7] | 5.28 | 17.00 | 7.58 |
| RCPR [21] | 6.18 | 17.26 | 8.35 |
| SDM [9] | 5.57 | 15.40 | 7.50 |
| LBF [8] | 4.95 | 11.98 | 6.32 |
| CFSS [10] | **4.73** | **9.98** | **5.76** |
| 3DDFA [14] | 6.15 | 10.59 | 7.01 |
| 3DDFA+SDM [14] | 5.53 | **9.56** | 6.31 |
| Ours (SDM+DFF) | **4.91** | 11.61 | **6.22** |

Figure 11: Examples of 2D and 3D facial alignment results on the AFLW2000 dataset [14] using our method, and the methods of Zhu et al. [14], Jourabloo & Liu [15], and Kazemi et al. [45].
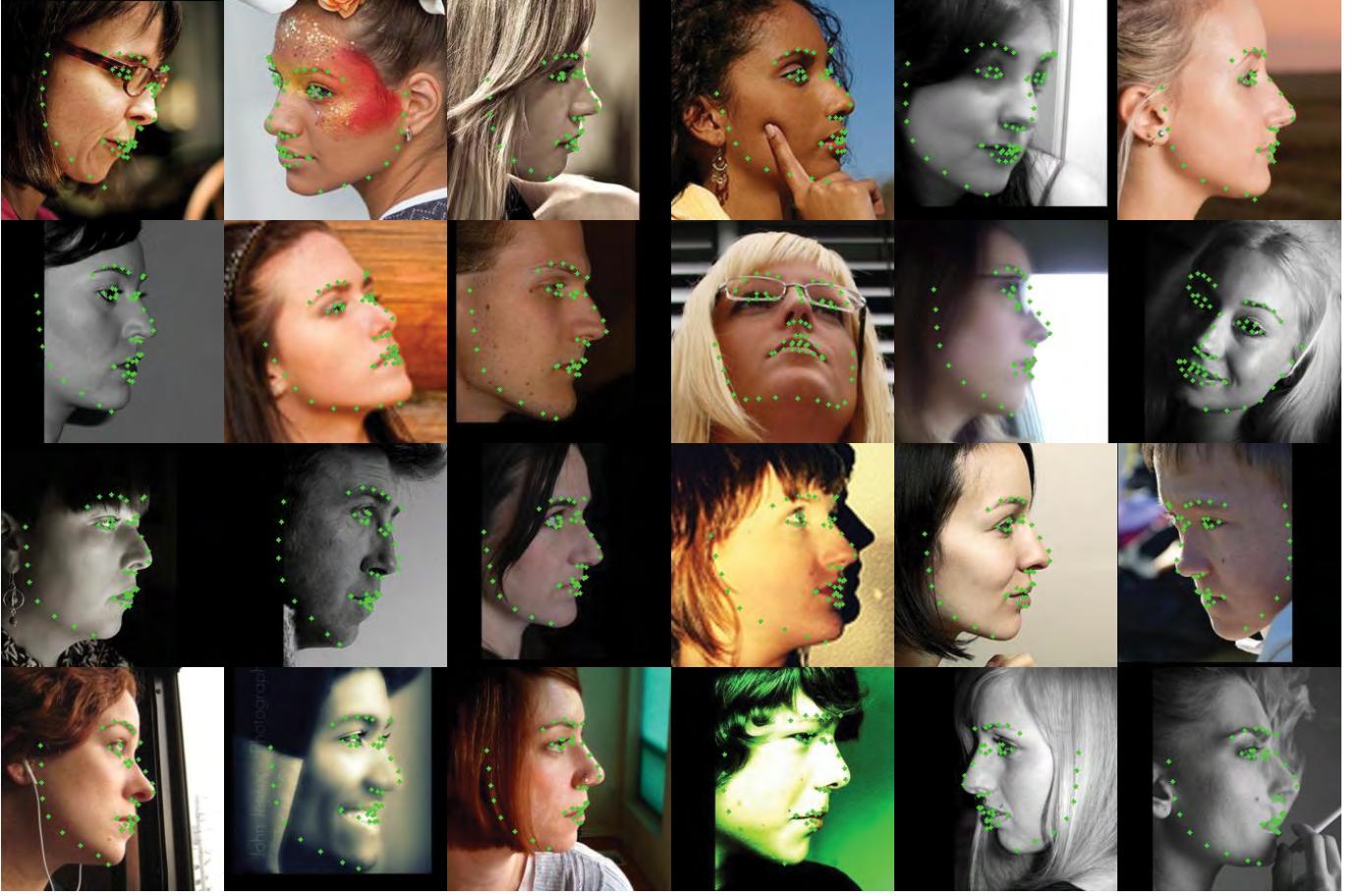
Figure 12: Examples of face alignment on large poses from the AFLW database. Only visible landmarks are showed.

improves the performance of SDM by simply replacing SIFT with DFF, achieving top performance for both the common and full test sets, and close to the top result for challenging test set.

The accompanying video[3] shows face feature tracking results using our method, on video clips with speech, poor lighting, large poses, and sports scenes. For comparison, we also show the results using [14], [15], [45]. Some parts of the video clips with fast face movements are slowed down for more clear comparison. Our results are consistently more robust and accurate for these challenging scenes.

### C. Computation Time

Our experiments are run on a PC with an Intel Core i7-4790 CPU at 4.0GHz, 8GB RAM, and a GTX 1070 GPU. It takes about 5ms to extract DFF for a $224 \times 224$ image, and about 8.8ms for each iteration of method described in III-C. All the results by our method are generated using three iterations. As we only need to extract DFF once for each image, the total running time is less than 35ms. Our method takes less computation time than the method in [14], while achieving more accurate alignment results.

## V. CONCLUSIONS

We present a deep learning based method to extract features from face images. Using a novel feature training method that utilizes the ground-truth correspondence between face images under different poses and expressions in the training set, the resulting deep face feature (DFF) has similar values for the same semantic pixels from different face images. As a result, the DFF captures global structural information of face images, and is more effective than general feature descriptors like SIFT on face related tasks such as matching and alignment. We propose a new face alignment method based on the DFF, which achieves state-of-the-art results for large-pose face alignment. In the future, it would be interesting to explore the use of DFF in other face related tasks.

## REFERENCES

[1] F. Liu, D. Zeng, Q. Zhao, and X. Liu, "Joint face alignment and 3d face reconstruction," in *European Conference on Computer Vision*. Springer, 2016, pp. 545–560.

[2] A. Wagner, J. Wright, A. Ganesh, Z. Zhou, H. Mobahi, and Y. Ma, "Toward a practical face recognition system: Robust alignment and illumination by sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 2, pp. 372–386, 2012.

[3] G. Tzimiropoulos and M. Pantic, "Optimization problems for fast AAM fitting in-the-wild," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 593–600.

[4] J. Saragih and R. Goecke, "A nonlinear discriminative approach to AAM fitting," in *Proceedings of the International Conference on Computer Vision*. IEEE, 2007, pp. 1–8.

[5] T. F. Cootes, G. J. Edwards, and C. J. Taylor, "Active appearance models," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 23, no. 6, pp. 681–685, 2001.

[6] I. Matthews and S. Baker, "Active appearance models revisited," *International journal of computer vision*, vol. 60, no. 2, pp. 135–164, 2004.

[7] X. Cao, Y. Wei, F. Wen, and J. Sun, "Face alignment by explicit shape regression," *International Journal of Computer Vision*, vol. 107, no. 2, pp. 177–190, 2014.

[8] S. Ren, X. Cao, Y. Wei, and J. Sun, "Face alignment at 3000 fps via regressing local binary features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1685–1692.

[9] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 532–539.

[10] S. Zhu, C. Li, C. Change Loy, and X. Tang, "Face alignment by coarse-to-fine shape searching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4998–5006.

[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, SIGGRAPH*, 1999, pp. 187–194.

[12] A. Jourabloo and X. Liu, "Pose-invariant 3D face alignment," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3694–3702.

[13] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.

[14] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li, "Face alignment across large poses: A 3d solution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 146–155.

[15] 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*, 2016, pp. 4188–4196.

[16] R. Yu, S. Saito, H. Li, D. Ceylan, and H. Li, "Learning dense facial correspondences in unconstrained images," *arXiv preprint arXiv:1709.00536*, 2017.

[17] T. F. Cootes, C. J. Taylor, and A. Lanitis, "Active shape models: Evaluation of a multi-resolution method for improving image search," in *Proc. British Machine Vision Conference*, 1994, pp. 327–338.

[18] D. Cristinacce and T. F. Cootes, "Boosted regression active shape models," in *Proc. British Machine Vision Conference*, 2007, pp. 880–889.

[19] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic, "Robust discriminative response map fitting with constrained local models," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 3444–3451.

[20] D. Cristinacce and T. Cootes, "Automatic feature localisation with constrained local models," *Pattern Recognition*, vol. 41, no. 10, pp. 3054–3067, 2008.

[21] X. P. Burgos-Artizzu, P. Perona, and P. Dollár, "Robust face landmark estimation under occlusion," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1513–1520.

[22] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3476–3483.

[23] J. Zhang, S. Shan, M. Kan, and X. Chen, "Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment," in *European Conference on Computer Vision*. Springer, 2014, pp. 1–16.

[24] E. Zhou, H. Fan, Z. Cao, Y. Jiang, and Q. Yin, "Extensive facial landmark localization with coarse-to-fine convolutional network cascade," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, pp. 386–391.

[25] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "Facial landmark detection by deep multi-task learning," in *European Conference on Computer Vision*. Springer, 2014, pp. 94–108.

[26] X. Yu, J. Huang, S. Zhang, W. Yan, and D. N. Metaxas, "Pose-free facial landmark fitting via optimized part mixtures and cascaded deformable shape model," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1944–1951.

[27] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 2879–2886.

[28] L. A. Jeni, J. F. Cohn, and T. Kanade, "Dense 3D face alignment from 2D videos in real-time," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, vol. 1. IEEE, 2015, pp. 1–8.

[29] V. Blanz and T. Vetter, "Face recognition based on fitting a 3D morphable model," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 25, no. 9, pp. 1063–1074, 2003.

[30] 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*. IEEE, 2005, pp. 986–993.

[31] C. Wu, S. Agarwal, B. Curless, and S. M. Seitz, "Multicore bundle adjustment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2011, pp. 3057–3064.

[32] C. Wu, "Towards linear-time incremental structure from motion," in *3DTV-Conference, 2013 International Conference on*. IEEE, 2013, pp. 127–134.

[33] M. Meytlis and L. Sirovich, "On the dimensionality of face space," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 7, pp. 1262–1267, 2007.

[34] J. P. Lewis, K. Anjyo, T. Rhee, M. Zhang, F. Pighin, and Z. Deng, "Practice and theory of blendshape facial models," in *Eurographics 2014 – State of the Art Reports*, S. Lefebvre and M. Spagnuolo, Eds. The Eurographics Association, 2014.

[35] 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*. IEEE, 2009, pp. 296–301.

[36] C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou, "Facewarehouse: A 3d facial expression database for visual computing," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 3, pp. 413–425, 2014.

[37] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic, "300 faces in-the-wild challenge: The first facial landmark localization challenge," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, pp. 397–403.

[38] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker, "Multi-pie," *Image and Vision Computing*, vol. 28, no. 5, pp. 807–813, 2010.

[39] Y. Guo, J. Zhang, J. Cai, B. Jiang, and J. Zheng, "3dfacenet: Real-time dense face reconstruction via synthesizing photo-realistic face images," *arXiv preprint arXiv:1708.00980*, 2017.

[40] S. Ren, X. Cao, Y. Wei, and J. Sun, "Face alignment via regressing local binary features," *IEEE Trans. Image Processing*, vol. 25, no. 3, pp. 1233–1245, 2016.

[41] L. Wei, Q. Huang, D. Ceylan, E. Vouga, and H. Li, "Dense human body correspondences using convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1544–1553.

[42] Q. Du, V. Faber, and M. Gunzburger, "Centroidal voronoi tessellations: Applications and algorithms," *SIAM Review*, vol. 41, no. 4, pp. 637–676, 1999.

[43] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI (3)*, ser. Lecture Notes in Computer Science, vol. 9351. Springer, 2015, pp. 234–241.

[44] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6738–6746.

[45] V. Kazemi and J. Sullivan, "One millisecond face alignment with an ensemble of regression trees," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1867–1874.