# Deep Multi-Center Learning for Face Alignment

Zhiwen Shao[1*], Hengliang Zhu[1], Xin Tan[1], Yangyang Hao[1], and Lizhuang Ma[2,1*]

[1]Department of Computer Science and Engineering, Shanghai Jiao Tong University, China

[2]School of Computer Science and Software Engineering, East China Normal University, China

{shaozhiwen, hengliang_zhu, tanxin2017, haoyangyang2014}@sjtu.edu.cn, ma-lz@cs.sjtu.edu.cn

***Abstract*—Facial landmarks are highly correlated with each other since a certain landmark can be estimated by its neighboring landmarks. Most of the existing deep learning methods only use one fully-connected layer called shape prediction layer to estimate the locations of facial landmarks. In this paper, we propose a novel deep learning framework named Multi-Center Learning with multiple shape prediction layers for face alignment. In particular, each shape prediction layer emphasizes on the detection of a certain cluster of semantically relevant landmarks respectively. Challenging landmarks are focused firstly, and each cluster of landmarks is further optimized respectively. Moreover, to reduce the model complexity, we propose a model assembling method to integrate multiple shape prediction layers into one shape prediction layer. Extensive experiments demonstrate that our method is effective for handling complex occlusions and appearance variations with real-time performance. The code for our method is available at https://github.com/ZhiwenShao/MCNet-Extension.**

***Index Terms*—Multi-Center Learning, Model Assembling, Face Alignment**

## I. INTRODUCTION

Face alignment refers to detecting facial landmarks such as eye centers, nose tip, and mouth corners. It is the preprocessor stage of many face analysis tasks like face animation [1], face beautification [2], and face recognition [3]. A robust and accurate face alignment is still challenging in unconstrained scenarios, owing to severe occlusions and large appearance variations. Most conventional methods [4], [5], [6], [7] only use low-level handcrafted features and are not based on the prevailing deep neural networks, which limits their capacity to represent highly complex faces.

Recently, several methods use deep networks to estimate shapes from input faces. Sun et al. [8], Zhou et al. [9], and Zhang et al. [10] employed cascaded deep networks to refine predicted shapes successively. Due to the use of multiple networks, these methods have high model complexity with complicated training processes. Taking this into account, Zhang et al. [11], [12] proposed a Tasks-Constrained Deep Convolutional Network (TCDCN), which uses only one deep network with excellent performance. However, it needs extra labels of facial attributes for training samples, which limits its universality.

Each facial landmark is not isolated but highly correlated with adjacent landmarks. As shown in Fig. 1(a), facial landmarks along the chin are all occluded, and landmarks around the mouth are partially occluded. Fig. 1(b) shows that landmarks on the right side of face are almost invisible. Therefore,



(a) Chin is occluded.　(b) Right contour is invisible.

Fig. 1. Examples of unconstrained face images with partial occlusions and large pose.

landmarks in the same local face region have similar properties including occlusion and visibility. It is observed that the nose can be localized roughly with the locations of eyes and mouth. There are also structural correlations among different facial parts. Motivated by this fact, facial landmarks are divided into several clusters based on their semantic relevance.

In this work[1], we propose a novel deep learning framework named Multi-Center Learning (MCL) to exploit the strong correlations among landmarks. In particular, our network uses multiple shape prediction layers to predict the locations of landmarks, and each shape prediction layer emphasizes on the detection of a certain cluster of landmarks respectively. By weighting the loss of each landmark, challenging landmarks are focused firstly, and each cluster of landmarks is further optimized respectively. Moreover, to decrease the model complexity, we propose a model assembling method to integrate multiple shape prediction layers into one shape prediction layer. The entire framework reinforces the learning process of each landmark with a low model complexity.

The main contributions of this study can be summarized as follows:

- We propose a novel multi-center learning framework for exploiting the strong correlations among landmarks.
- We propose a model assembling method which ensures a low model complexity.
- Extensive experiments demonstrate that our method is effective for handling complex occlusions and appearance variations with real-time performance.

---

[1]This is an extended version of [13] with two improvements. The shape prediction layer is replaced from the fully-connected layer to the Global Average Pooling layer [14], which has a stronger feature learning ability. To exploit the correlations among landmarks more completely, challenging landmarks are focused firstly before each cluster of landmarks is respectively optimized.

The remainder of this paper is structured as below. We discuss related works in the next section. In Section III, we illuminate the structure of our network and the learning algorithm. Extensive experiments are carried out in Section IV. Section V concludes this work.

## II. RELATED WORK

We review researches from three aspects related to our method: conventional face alignment, unconstrained face alignment, face alignment via deep learning.

### A. Conventional Face Alignment

Conventional face alignment methods can be classified as two categories: template fitting and regression-based.

Template fitting methods match faces by constructing shape templates. Cootes et al. [15] proposed a typical template fitting method named Active Appearance Model (AAM), which minimizes the texture residual to estimate the shape. Asthana et al. [16] used regression techniques to learn functions from response maps to shapes, in which the response map has stronger robustness and generalization ability than texture based features of AAM. Pedersoli et al. [17] developed the mixture of trees of parts method by extending the mixtures from trees to graphs, and learned a deformable detector to align its parts to faces. However, these templates are not complete enough to cover complex variations, which are difficult to be generalized to unseen faces.

Regression-based methods predict the locations of facial landmarks by learning a regression function from face features to shapes. Cao et al. [4] proposed an Explicit Shape Regression (ESR) method to predict the shape increment with pixel-difference features. Xiong et al. [5] proposed a Supervised Descent Method (SDM) to detect landmarks by solving the nonlinear least squares problem, with Scale-Invariant Feature Transform (SIFT) [18] features and linear regressors being applied. Ren et al. [7] used a locality principle to extract a set of Local Binary Features (LBF), in which a linear regression is utilized for localizing landmarks. Lee et al. [19] employs Cascade Gaussian Process Regression Trees (cGPRT) with shape-indexed difference of Gaussian features to achieve face alignment. It has a better generalization ability than cascade regression trees, and shows strong robustness against geometric variations of faces. Most of these methods give an initial shape and refine the shape in an iterative manner, where the final solutions are prone to getting trapped in a local optimum with a poor initialization. In contrast, our method uses a deep neural network to regress from raw face patches to the locations of landmarks.

### B. Unconstrained Face Alignment

Large pose variations and severe occlusions are major challenges in unconstrained environments. Unconstrained face alignment methods are based on 3D models or deal with occlusions explicitly.

Many methods utilize 3D shape models to solve large-pose face alignment. Nair et al. [20] refined the fit of a 3D point distribution model to perform landmark detection. Yu et al. [21] used a cascaded deformable shape model to detect landmarks of large-pose faces. Cao et al. [1] employed a displaced dynamic expression regression to estimate the 3D face shape and 2D facial landmarks. The predicted 2D landmarks are used to adjust the model parameters to better fit the current user. Jeni et al. [22] proposed a 3D cascade regression method to implement 3D face alignment, which can maintain the pose invariance of facial landmarks within the range of around 60 degrees.

There are several occlusion-free face alignment methods. Burgos-Artizzu et al. [6] developed a Robust Cascaded Pose Regression (RCPR) method to detect occlusions explicitly, and uses shape-indexed features to regress the shape increment. Yu et al. [23] utilizes a Bayesian model to merge the estimation results from multiple regressors, in which each regressor is trained to localize facial landmarks with a specific pre-defined facial part being occluded. Wu et al. [24] proposed a Robust Facial Landmark Detection (RFLD) method, which uses a robust cascaded regressor to handle complex occlusions and large head poses. To improve the performance of occlusion estimation, landmark visibility probabilities are estimated with an explicit occlusion constraint. Different from these methods, our method is not based on 3D models and does not process occlusions explicitly.

### C. Face Alignment via Deep Learning

Deep learning methods can be divided into two classes: single network based and multiple networks based.

Sun et al. [8] estimated the locations of 5 facial landmarks using Cascaded Convolutional Neural Networks (Cascaded CNN), in which each level computes averaged estimated shape and the shape is refined level by level. Zhou et al. [9] used multi-level deep networks to detect facial landmarks from coarse to fine. Similarly, Zhang et al. [10] proposed Coarse-to-Fine Auto-encoder Networks (CFAN). These methods all use multi-stage deep networks to localize landmarks in a coarse-to-fine manner. Instead of using cascaded networks, Honari et al. [25] proposed Recombinator Networks (RecNet) for learning coarse-to-fine feature aggregation with multi-scale input maps, where each branch extracts features based on current maps and the feature maps of coarser branches.

A few methods employ a single network to solve the face alignment problem. Shao et al. [26] proposed a Coarse-to-Fine Training (CFT) method to learn the mapping from input face patches to estimated shapes, which searches the solutions smoothly by adjusting the relative weights between principal landmarks and elaborate landmarks. Zhang et al. [11], [12] used the TCDCN with auxiliary facial attribute recognition to predict correlative facial properties like expression and pose, which improves the performance of face alignment. Xiao et al. [27] proposed a Recurrent Attentive-Refinement (RAR) network for face alignment under unconstrained conditions, where shape-indexed deep features and temporal information are taken as inputs and shape predictions are recurrently revised. Compared to these methods, our method uses only one network and is independent of additional facial attributes.
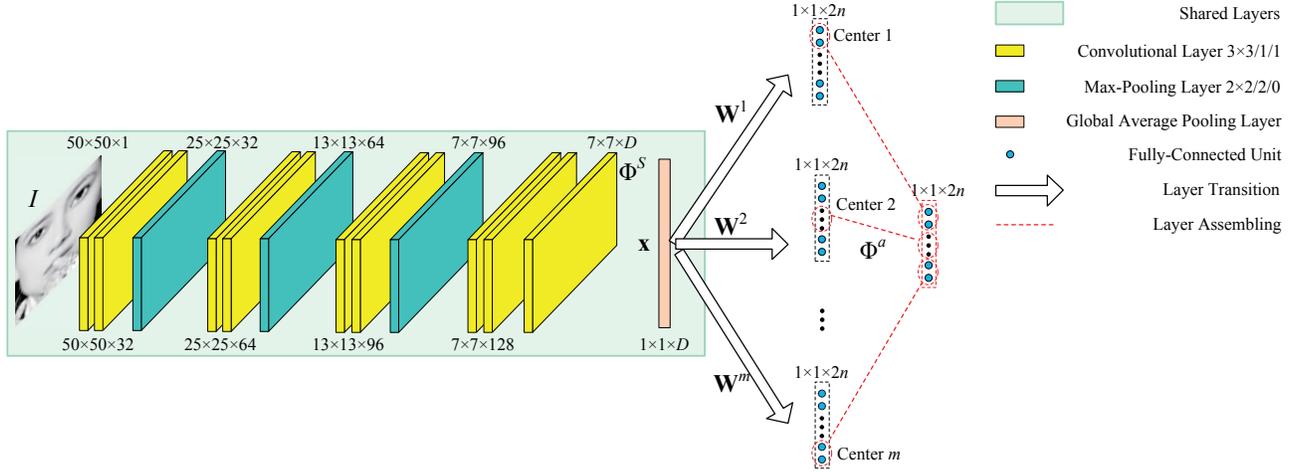
Fig. 2. Architecture of our network MCL. The expression $k_1 \times k_2 \times k_3$ attached to each layer denotes the height, width, and channel respectively. Every two convolutional layers possess the same expression. The expression $k_4 \times k_5/k_6/k_7$ denotes the height, width, stride, and padding of filters respectively. The same type of layers use the identical filters.

## III. MULTI-CENTER LEARNING FOR FACE ALIGNMENT

### A. Network Architecture

The architecture of our network MCL is illustrated in Fig. 2. MCL contains three max-pooling layers, each of which follows a stack of two convolutional layers proposed by VGGNet [28]. In the fourth stack of convolutional layers, we use a convolutional layer with $D$ feature maps above two convolutional layers. We perform Batch Normalization (BN) [29] and Rectified Linear Unit (ReLU) [30] after each convolution to accelerate the convergence of our network. Most of the existing deep learning methods such as TCDCN [11], [12] use the fully-connected layer to extract features, which is apt to overfit and hamper the generalization ability of the network. To sidestep these problems, we operate Global Average Pooling [14] on the last convolutional layer to extract a high-level feature representation $\mathbf{x}$, which computes the average of each feature map. With this improvement, our MCL acquires a higher representation power with fewer parameters.

Face alignment can be regarded as a nonlinear regression problem, which transforms appearance to shape. A transformation $\Phi^S(\cdot)$ is used for modeling this highly nonlinear function, which extracts the feature $\mathbf{x}$ from the input face image $I$, formulated as

$$\mathbf{x} = \Phi^S(I), \tag{1}$$

where $\mathbf{x} = (x_0, x_1, \cdots, x_D)^T \in \mathbb{R}^{(D+1) \times 1}$, $x_0 = 1$ corresponds to the bias, and $\Phi^S(\cdot)$ is a composite function of operations including convolution, BN, ReLU, and pooling.

Traditionally, only one shape prediction layer is used, which limits the performance. In contrast, our MCL uses multiple shape prediction layers, each of which emphasizes on the detection of a certain cluster of landmarks. The first several layers are shared by multiple shape prediction layers, which are called shared layers forming the composite function $\Phi^S(\cdot)$. For the $i$-th shape prediction layer, $i = 1, \cdots, m$, a weight matrix $\mathbf{W}^i = (\mathbf{w}_1^i, \mathbf{w}_2^i, \cdots, \mathbf{w}_{2n}^i) \in \mathbb{R}^{(D+1) \times 2n}$ is used to connect the feature $\mathbf{x}$, where $m$ and $n$ are the number of shape prediction layers and landmarks, respectively. The reason why we train each shape prediction layer to predict $n$ landmarks instead of one cluster of landmarks is that different facial parts have correlations, as shown in Fig. 1.

To decrease the model complexity, we use a model assembling function $\Phi^a(\cdot)$ to integrate multiple shape prediction layers into one shape prediction layer, which is formulated as

$$\mathbf{W}^a = \Phi^a(\mathbf{W}^1, \cdots, \mathbf{W}^m), \tag{2}$$

where $\mathbf{W}^a = (\mathbf{w}_1^a, \mathbf{w}_2^a, \cdots, \mathbf{w}_{2n}^a) \in \mathbb{R}^{(D+1) \times 2n}$ is the assembled weight matrix. Specifically, $\mathbf{w}_{2j-1}^a = \mathbf{w}_{2j-1}^i$, $\mathbf{w}_{2j}^a = \mathbf{w}_{2j}^i$, $j \in P^i$, $i = 1, \cdots, m$, where $P^i$ is the $i$-th cluster of indexes of landmarks. The final prediction $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_{2n})$ is defined as

$$\hat{\mathbf{y}} = \mathbf{W}^{aT}\mathbf{x}, \tag{3}$$

where $\hat{y}_{2j-1}$ and $\hat{y}_{2j}$ denote the predicted x-coordinate and y-coordinate of the $j$-th landmark respectively.

Compared to other typical convolutional networks like VG-GNet [28], GoogLe-Net [31], and ResNet [32], our network MCL is substantially smaller and shallower. We believe that such a concise structure is efficient for estimating the locations of facial landmarks. Firstly, face alignment aims to regress coordinates of fewer than 100 facial landmarks generally, which demands much lower model complexity than visual recognition problems with more than $1,000$ classes. Secondly, a very deep network may fail to work well for landmark detection owing to the reduction of spatial information layer by layer. Other visual localization tasks, like face detection, usually use multiple cascaded shallow networks rather than a single very deep network. Finally, common face alignment benchmarks only contain thousands of training images. A simple network is not easy to overfit given a small amount of raw training data.

### B. Learning Algorithm

The overview of our learning algorithm is shown in Algorithm 1. $\Omega^t$ and $\Omega^v$ are the training set and the validation

**Algorithm 1** Multi-Center Learning Algorithm.

**Input:** A network MCL, $\Omega^t$, $\Omega^v$, initialized $\Theta$.
**Output:** $\Theta$.
 1: Pre-train shared layers and one shape prediction layer until convergence;
 2: Fix the parameters of the first six convolutional layers and fine-tune subsequent layers until convergence;
 3: Fine-tune all the layers until convergence;
 4: **for** $i = 1$ to $m$ **do**
 5:   Fix $\Theta^S$ and fine-tune the $i$-th shape prediction layer until convergence;
 6: **end for**
 7: $\Theta = \Theta^S \cup \mathbf{W}^a$;
 8: Return $\Theta$.

set respectively. $\Theta$ is the set of parameters including weights and biases of our network MCL, which is updated using Mini-Batch Stochastic Gradient Descent (SGD) [33] at each iteration. The face alignment loss is defined as

$$E = \sum_{j=1}^{n} u_j[(y_{2j-1} - \hat{y}_{2j-1})^2 + (y_{2j} - \hat{y}_{2j})^2]/(2d^2), \quad (4)$$

where $u_j$ is the weight of the $j$-th landmark, $y_{2j-1}$ and $y_{2j}$ denote the ground-truth x-coordinate and y-coordinate of the $j$-th landmark respectively, and $d$ is the ground truth inter-ocular distance between the eye centers.

Inter-ocular distance normalization provides fair comparisons among faces with different size, and reduces the magnitude of loss to speed up the learning process. During training, a too high learning rate may cause the missing of optimum so far as to the divergence of network, and a too low learning rate may lead to falling into a local optimum. We employ a low initial learning rate to avoid the divergence, and increase the learning rate when the loss is reduced significantly and continue the training procedure.

*1) Pre-Training and Weighting Fine-Tuning:* In Step 1, a *basic model (BM)* with one shape prediction layer is pre-trained to learn a good initial solution. In Eq. 4, $u_j = 1$ for all $j$. The average alignment error of each landmark of BM on $\Omega^v$ are $\epsilon_1^b, \cdots, \epsilon_n^b$ respectively, which are averaged over all the images. The landmarks with larger errors than remaining landmarks are treated as challenging landmarks.

In Steps 2 and 3, we focus on the detection of challenging landmarks by assigning them larger weights. The weight of the $j$-th landmark is proportional to its alignment error as

$$u_j = n\epsilon_j^b / \sum_{j=1}^{n} \epsilon_j^b. \quad (5)$$

Instead of fine-tuning all the layers from BM directly, we use two steps to search the solution smoothly. Step 2 searches the solution without deviating from BM overly. Step 3 searches the solution within a larger range on the basis of the previous step. This stage is named weighting fine-tuning, which learns a *weighting model (WM)* with higher localization accuracy of challenging landmarks.

*2) Multi-Center Fine-Tuning and Model Assembling:* The face is partitioned into seven parts according to its semantic structure: left eye, right eye, nose, mouth, left contour, chin, and right contour. As shown in Fig. 3, different labeling patterns of 5, 29, and 68 facial landmarks are partitioned into 4, 5, and 7 clusters respectively. For the $i$-th shape prediction layer, the $i$-th cluster of landmarks are treated as the optimized center, and the set of indexes of remaining landmarks is denoted as $Q^i$.



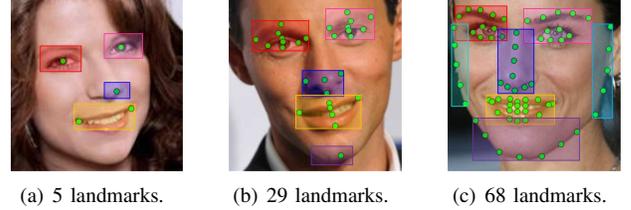(a) 5 landmarks.     (b) 29 landmarks.     (c) 68 landmarks.

Fig. 3. Partitions of facial landmarks for different labeling patterns.

From Steps 4 to 6, the parameters of shared layers $\Theta^S$ are fixed, and each shape prediction layer is initialized with the parameters of the shape prediction layer of WM. When fine-tuning the $i$-th shape prediction layer, the weights of landmarks in $P^i$ and $Q^i$ are defined as

$$u_{P^i} = \alpha u_{Q^i}, \quad (6)$$

where $\alpha \gg 1$ is a coefficient to make the $i$-th shape prediction layer emphasize on the detection of the $i$-th cluster of landmarks. The constraint between $u_{P^i}$ and $u_{Q^i}$ is formulated as

$$u_{P^i}|P^i| + u_{Q^i}(n - |P^i|) = n, \quad (7)$$

where $|\cdot|$ refers to the number of elements in a cluster. With Eqs. 6 and 7, the solved weights are formulated as

$$u_{P^i} = \alpha n/[(\alpha - 1)|P^i| + n],$$
$$u_{Q^i} = n/[(\alpha - 1)|P^i| + n]. \quad (8)$$

The average alignment error of each landmark of WM on $\Omega^v$ are $\epsilon_1^w, \cdots, \epsilon_n^w$ respectively. Similar to Eq. 5, the weight of the $j$-th landmark is

$$u_j = \begin{cases} u_{P^i}|P^i| \cdot \epsilon_j^w / \sum_{j \in P^i} \epsilon_j^w, & j \in P^i, \\ u_{Q^i}(n - |P^i|) \cdot \epsilon_j^w / \sum_{j \in Q^i} \epsilon_j^w, & j \in Q^i. \end{cases} \quad (9)$$

Although the landmarks in $P^i$ are mainly optimized, remaining landmarks are still considered with very small weights rather than zero. This is beneficial for utilizing implicit structural correlations of different facial parts and searching the solutions smoothly. This stage is called multi-center fine-tuning which learns multiple shape prediction layers.

In Step 7, multiple shape prediction layers are assembled into one shape prediction layer by Eq. 2. With this model assembling stage, our method learns an *assembling model (AM)*. There is no increase of model complexity in the assembling process, so AM has a low computational cost. It improves the detection precision of each facial landmark by integrating the advantage of each shape prediction layer.

*3) Analysis of Model Learning:* To investigate the influence from the weights of landmarks on learning procedure, we calculate the derivative of Eq. 4 with respect to $\hat{y}_k$:

$$\frac{\partial E}{\partial \hat{y}_k} = u_j(\hat{y}_k - y_k)/d^2, \tag{10}$$

where $k \in \{2j-1, 2j\}$, $j = 1, \cdots, n$. During the learning process, the assembled weight matrix $\mathbf{W}^a$ in Eq. 3 is updated by SGD. Specifically, $\mathbf{w}_k^a = \mathbf{w}_k^a - \eta \frac{\partial E}{\partial \mathbf{w}_k^a} = \mathbf{w}_k^a - \eta \frac{\partial E}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \mathbf{w}_k^a} = \mathbf{w}_k^a - \eta \frac{\partial E}{\partial \hat{y}_k} \mathbf{x}$. In summary, $\mathbf{W}^a$ is updated as

$$\mathbf{w}_k^a = \mathbf{w}_k^a - \eta u_j(\hat{y}_k - y_k)\mathbf{x}/d^2, \tag{11}$$

where $\eta$ is the learning rate. If the $j$-th landmark is given a larger weight, its corresponding parameters will be updated with a larger step towards the optimal solution. Therefore, weighting the loss of each landmark ensures that the landmarks with larger weights are mainly optimized. Our method first uses the weighting fine-tuning stage to optimize challenging landmarks, and further uses the multi-center fine-tuning stage to optimize each cluster of landmarks respectively.

## IV. EXPERIMENTS

### A. Datasets and Settings

*1) Datasets:* There are three challenging benchmarks AFLW [34], COFW [6], and IBUG [35], which are used for evaluating face alignment with severe occlusion and large variations of pose, expression, and illumination. The provided face bounding boxes are employed to crop face patches during testing.

- **AFLW [34]** contains $25,993$ faces under real-world conditions gathered from Flickr. Compared with other datasets like MUCT [36] and LFPW [37], AFLW exhibits larger pose variations and extreme partial occlusions. Following the settings of [11], [25], $2,995$ images are used for testing, and $10,000$ images annotated with 5 landmarks are used for training, which includes $4,151$ LFW [38] images and $5,849$ web images.
- **COFW [6]** is an occluded face dataset in the wild, in which the faces are designed with severe occlusions using accessories and interactions with objects. It contains $1,007$ images annotated with 29 landmarks. The training set includes 845 LFPW faces and 500 COFW faces, and the testing set includes remaining 507 COFW faces.
- **IBUG [35]** contains 135 testing images which present large variations in pose, expression, illumination, and occlusion. The training set consists of AFW [39], the training set of LFPW, and the training set of Helen [40], which are from 300-W [35] with $3,148$ images labeled with 68 landmarks.

*2) Implementation Details:* We enhance the diversity of raw training data on account of their limited variation patterns, using five steps: rotation, uniform scaling, translation, horizontal flip, and JPEG compression. In particular, for each training face, we firstly perform multiple rotations, and attain a tight face bounding box covering the ground truth locations of landmarks of each rotated result respectively. Uniform scaling and translation with different extents on face bounding boxes

are further conducted, in which each newly generated face bounding box is used to crop the face. Finally training samples are augmented through horizontal flip and JPEG compression. It is beneficial for avoiding overfitting and improving the robustness of learned models by covering various patterns.

We train our MCL using an open source deep learning framework Caffe [41]. The input face patch is a $50 \times 50$ grayscale image, and each pixel value is normalized to $[-1, 1)$ by subtracting 128 and multiplying $0.0078125$. A more complex model is needed for a labeling pattern with more facial landmarks, so $D$ is set to be $512/512/1,024$ for $5/29/68$ facial landmarks. The type of solver is SGD with a mini-batch size of 64, a momentum of $0.9$, and a weight decay of $0.0005$. The maximum learning iterations of pre-training and each fine-tuning step are $18 \times 10^4$ and $6 \times 10^4$ respectively, and the initial learning rates of pre-training and each fine-tuning step are $0.02$ and $0.001$ respectively. Note that the initial learning rate of fine-tuning should be low to preserve some representational structures learned in the pre-training stage and avoid missing good intermediate solutions. The learning rate is multiplied by a factor of $0.3$ at every $3 \times 10^4$ iterations, and the remaining parameter $\alpha$ is set to be 125.

*3) Evaluation Metric:* Similar to previous methods [4], [8], [12], we report the inter-ocular distance normalized mean error, and treat the mean error larger than $10\%$ as a failure. To conduct a more comprehensive comparison, the cumulative errors distribution (CED) curves are plotted. To measure the time efficiency, the average running speed (Frame per Second, FPS) on a single core i5-6200U 2.3GHz CPU is also reported. A single image is fed into the model at a time when computing the speed. In other words, we evaluate methods on four popular metrics: mean error (%), failure rate (%), CED curves, and average running speed. In the next sections, % in all the results are omitted for simplicity.

### B. Comparison with State-of-the-Art Methods

We compare our work MCL against state-of-the-art methods including ESR [4], SDM [5], Cascaded CNN [8], RCPR [6], CFAN [10], LBF [7], cGPRT [19], CFSS [42], TCDCN [11], [12], ALR [43], CFT [26], RFLD [24], RecNet [25], RAR [27], and FLD+PDE [44]. All the methods are evaluated on testing images using the face bounding boxes provided by benchmarks. In addition to given training images, TCDCN uses outside training data labeled with facial attributes. RAR augments training images with occlusions incurred by outside natural objects like sunglasses, phones, and hands. FLD+PDE performs facial landmark detection, pose and deformation estimation simultaneously, in which the training data of pose and deformation estimation are used. Other methods including our MCL only utilize given training images from the benchmarks.

Table I reports the results of our method and previous works on three benchmarks. Our method MCL outperforms most of the state-of-the-art methods, especially on AFLW dataset where a relative error reduction of $3.93\%$ is achieved compared to RecNet. Cascaded CNN estimates the location of each

---

[2]The result is acquired by running the code at https://github.com/seetaface/SeetaFaceEngine/tree/master/FaceAlignment.

(a) Results of Cascaded CNN, ALR, and MCL on AFLW.
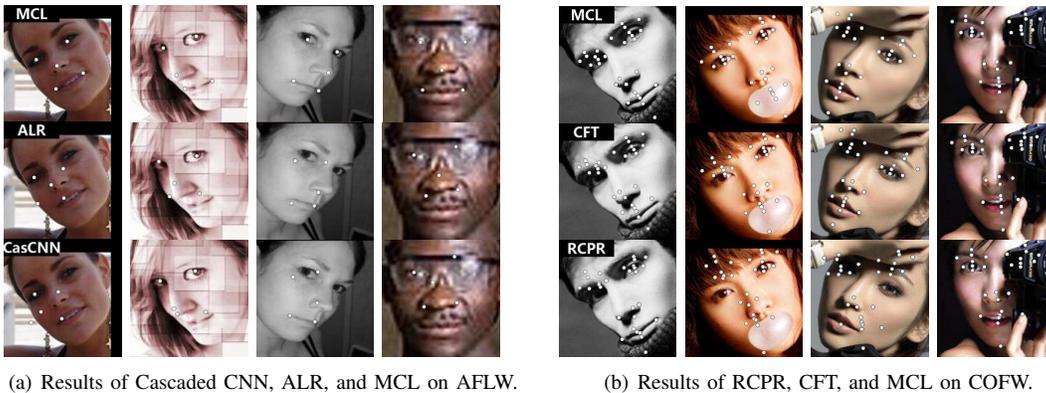
(b) Results of RCPR, CFT, and MCL on COFW.

Fig. 4. Face alignment results of state-of-the-art methods and our method MCL on challenging faces.

TABLE I

COMPARISON OF RESULTS OF MEAN ERROR WITH STATE-OF-THE-ART METHODS. SEVERAL METHODS DID NOT SHARE THEIR RESULTS OR CODE ON SOME BENCHMARKS, SO WE USE RESULTS MARKED WITH "*" FROM [11], [12].

| Method | AFLW 5 landmarks | COFW 29 landmarks | IBUG 68 landmarks |
|---|---|---|---|
| ESR [4] | 12.4* | 11.2* | 17.00* |
| SDM [5] | 8.5* | 11.14* | 15.40* |
| Cascaded CNN [8] | 8.72 | - | - |
| RCPR [6] | 11.6* | 8.5 | 17.26* |
| CFAN [10] | 7.83$^2$ | - | 16.78* |
| LBF [7] | - | - | 11.98 |
| cGPRT [19] | - | - | 11.03 |
| CFSS [42] | - | - | 9.98 |
| TCDCN [11], [12] | 8.0 | 8.05 | 8.60 |
| ALR [43] | 7.42 | - | - |
| CFT [26] | - | 6.33 | 10.06 |
| RFLD [24] | - | **5.93** | - |
| RecNet [25] | 5.60 | - | 8.44 |
| RAR [27] | 7.23 | 6.03 | **8.35** |
| FLD+PDE [44] | - | 6.40 | - |
| **MCL** | **5.38** | 6.00 | 8.51 |

landmark separately in the second and third level, and every two networks are used to detect one landmark. It is difficult to be extended to dense landmarks owing to the explosion of the number of networks. TCDCN relies on outside training data for auxiliary facial attribute recognition, which limits the universality. It can be seen that MCL outperforms Cascaded CNN and TCDCN on all the benchmarks. Moreover, MCL is robust to occlusions with the performance on par with RFLD, benefiting from utilizing semantical correlations among different landmarks. RecNet and RAR show significant results, but their models are very complex with high computational costs.

We compare with other methods on several challenging images from AFLW and COFW respectively in Fig. 4. Our method MCL indicates higher accuracy in the details than previous works. More examples on challenging IBUG are presented in Fig. 5. MCL demonstrates a superior capability of handling severe occlusions and complex variations of pose, expression, illumination. The CED curves of MCL and several state-of-the-art methods are shown in Fig. 6. It is observed that MCL achieves competitive performance on all three

benchmarks.

TABLE II

AVERAGE RUNNING SPEED OF DEEP LEARNING METHODS. THE TIME OF THE FACE DETECTION IS EXCLUDED.

| Method | Speed (FPS) | Platform |
|---|---|---|
| CFAN [10] | 43 | i7-3770 3.4 GHz CPU |
| TCDCN [12] | 50 | i5-6200U 2.3GHz CPU |
| CFT [26] | 31 | i5-6200U 2.3GHz CPU |
| RAR [27] | 4 | Titan-Z GPU |
| **MCL** | **57** | i5-6200U 2.3GHz CPU |

The average running speed of deep learning methods for detecting 68 facial landmarks are presented in Table II. Except for the methods tested on the i5-6200U 2.3GHz CPU, other methods are reported with the results in the original papers. Since CFAN utilizes multiple networks, it costs more running time. RAR achieves only 4 FPS on a Titan-Z GPU, which cannot be applied to practical scenarios. Both TCDCN and our method MCL are based on only one network, so they show higher speed. Our method only takes 17.5 ms per face on a single core i5-6200U 2.3GHz CPU. This profits from low model complexity and computational costs of our network. It can be concluded that our method is able to be extended to real-time facial landmark tracking in unconstrained environments.

TABLE III

RESULTS OF MEAN ERROR OF PRE-BM AND BM ON THREE BENCHMARKS.

| Method | AFLW 5 landmarks | COFW 29 landmarks | IBUG 68 landmarks |
|---|---|---|---|
| pre-BM [13] | **5.61** | 6.40 | 9.23 |
| **BM** | 5.67 | **6.25** | **8.89** |

### C. Ablation Study

*1) Global Average Pooling vs. Full Connection:* Based on the previous version of our work [13], the last max-pooling layer and the $D$-dimensional fully-connected layer are replaced with a convolutional layer and a Global Average Pooling layer[14]. The results of the mean error of BM and the previous version (pre-BM) [13] are shown in Table III. It can be seen that BM performs better on IBUG and

Fig. 5. Examples of LBF, CFSS, and MCL on challenging IBUG.



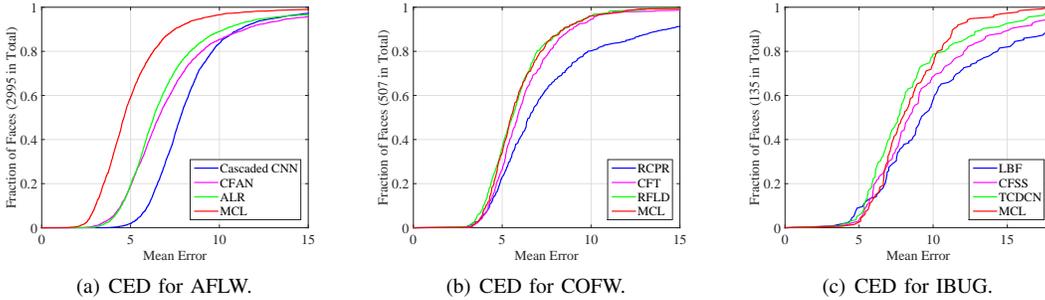(a) CED for AFLW.  (b) CED for COFW.  (c) CED for IBUG.

Fig. 6. Comparison of CED curves with previous methods on three benchmarks.

COFW but worse on AFLW than pre-BM. It demonstrates that Global Average Pooling is more advantageous for more complex problems with more facial landmarks. There are higher requirements for learned features when localizing more facial landmarks. For simple problems especially for localizing 5 landmarks of AFLW, a plain network with full connection is more prone to being trained.

The difference between pre-BM and BM is the structure of learning the feature $\mathbf{x}$. The number of parameters for this part of pre-BM and BM are $(4 \times 4 \times 128 + 1)D = 2,049D$ and $(3 \times 3 \times 128 + 1)D + 2D + 2D = 1,157D$ respectively, where the three terms for BM correspond to the convolution, the expectation and variance of BN [29], and the scaling and shifting of BN. Therefore, BM has a stronger feature learning ability with fewer parameters than pre-BM.
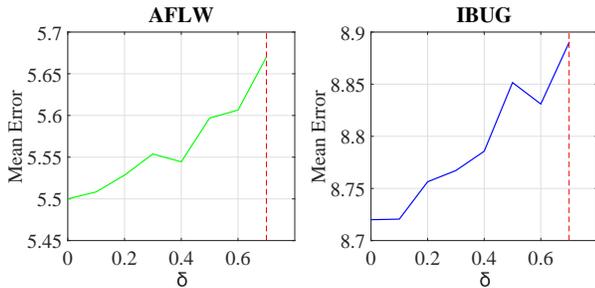


Fig. 7. Mean error of WM on AFLW and IBUG with different $\delta$.

*2) Robustness of Weighting:* To verify the robustness of weighting, random perturbations are added to the weights of landmarks. In particular, we plus a perturbation $\delta$ to the

weight of each of random $\lfloor n/2 \rfloor$ landmarks and minus $\delta$ to the weight of each of remaining $n - \lfloor n/2 \rfloor$ landmarks, where $\lfloor \cdot \rfloor$ refers to rounding down to the nearest integer. Fig. 7 shows the variations of mean error of WM with the increase of $\delta$. When $\delta$ is $0.4$, WM can still achieves good performance. Therefore, weighting the loss of each landmark is robust to random perturbations. Even if different weights are obtained, the results will not be affected as long as the relative sizes of weights are identical.

TABLE IV
RESULTS OF MEAN ERROR OF LANDMARKS OF EACH CLUSTER ON IBUG.

| Cluster | WM | Left Eye Model | Right Eye Model |
|---|---|---|---|
| Left Eye | 8.09 | **7.92** | 8.10 |
| Right Eye | 7.73 | 7.55 | **7.30** |
| Nose | 6.19 | 6.42 | 6.59 |
| Mouth | 6.92 | 6.80 | 7.08 |
| Left Contour | 12.66 | 12.83 | 12.74 |
| Chin | 13.55 | 13.50 | 13.45 |
| Right Contour | 13.38 | 13.47 | 13.45 |

*3) Analysis of Shape Prediction Layers:* Our method learns each shape prediction layer respectively with a certain cluster of landmarks being emphasized. The results of WM and two shape prediction layers with respect to the left eye and the right eye on IBUG benchmark are shown in Table IV. Compared to WM, the left eye model and the right eye model both reduce the alignment errors of their corresponding clusters. As a result, the assembled AM can improve the detection accuracy of landmarks of the left eye and the right eye on the basis of WM.

Note that the two models also improve the localization

precision of other clusters. Taking the left eye model as an example, it additionally reduces the errors of landmarks of right eye, mouth, and chin, which is due to the correlations among different facial parts. Moreover, for the right eye cluster, the right eye model improves the accuracy more significantly than the left eye model. It can be concluded that each shape prediction layer emphasizes on the corresponding cluster respectively.

| Method | Weighting Fine-Tuning | Multi-Center Fine-Tuning | COFW | IBUG |
|---|---|---|---|---|
| Simplified AM | | $\checkmark$ | 6.08 | 8.67 |
| **AM** | $\checkmark$ | $\checkmark$ | **6.00** | **8.51** |
| Weighting Simplified AM | $\checkmark$ | $\checkmark$ | 6.05 | 8.67 |

*4) Integration of Weighting Fine-Tuning and Multi-Center Fine-Tuning:* Here we validate the effectiveness of weighting fine-tuning by removing the weighting fine-tuning stage to learn a *Simplified AM* from BM. Table V presents the results of mean error of Simplified AM and AM respectively on COFW and IBUG. Note that Simplified AM has already acquired good results, which verifies the effectiveness of the multi-center fine-tuning stage. The accuracy of AM is superior to that of Simplified AM especially on challenging IBUG, which is attributed to the integration of two stages. A *Weighting Simplified AM* from Simplified AM using the weighting fine-tuning stage is also learned, whose results are shown in Table V. It can be seen that Weighting Simplified AM improves slightly on COFW but fails to search a better solution on IBUG. Therefore, we choose to use the multi-center fine-tuning stage after the weighting fine-tuning stage.

| Method | AFLW | | COFW | | IBUG | |
|---|---|---|---|---|---|---|
| | Error | Failure | Error | Failure | Error | Failure |
| BM | 5.67 | 4.43 | 6.25 | 5.13 | 8.89 | 27.43 |
| WM | 5.50 | 3.84 | 6.11 | 4.54 | 8.72 | 26.80 |
| **AM** | **5.38** | **3.47** | **6.00** | **3.94** | **8.51** | **25.93** |

*5) Discussion of All Stages:* Table VI summarizes the results of mean error and failure rate of BM, WM, and AM. It can be observed that AM has higher accuracy and stronger robustness than BM and WM. Fig. 8 depicts the enhancement from WM to AM for several examples of COFW. The localization accuracy of facial landmarks from each cluster is improved in the details. It is because each shape prediction layer increases the detection precision of corresponding cluster respectively.

### D. MCL for Partially Occluded Faces

The correlations among different facial parts are very useful for face alignment especially for partially occluded faces. To
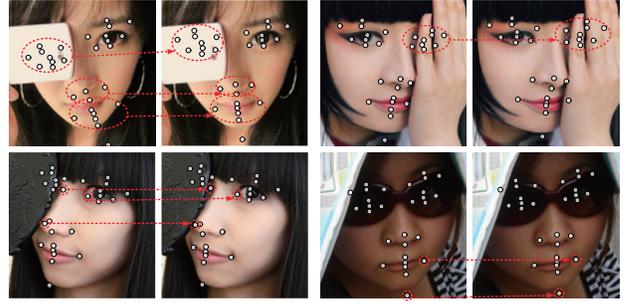


Fig. 8. Examples of improvement for different facial landmarks from WM to AM on COFW dataset.

investigate the influence of occlusions, we directly use trained WM and AM without any additional processing for partially occluded faces. Randomly 30% testing faces from COFW are processed with left eyes being occluded, where the tight bounding box covering landmarks of left eye is filled with gray color, as shown in Fig. 9.



Fig. 9. Example faces from COFW with left eyes occluded.

Table VII shows the mean error results for the left eye cluster and other clusters of WM and AM on COFW benchmark, where "with (w/) occlusion (occlu.)" denotes that left eyes of the testing faces are processed with handcrafted occlusions as illustrated in Fig. 9, and "without (w/o) occlu." denotes that the testing faces are kept unchanged. Note that our method does not process occlusions explicitly, in which the training data is not performed handcrafted occlusions. After processing testing faces with occlusions, the mean error results of both WM and AM increase. Besides the results of landmarks from the left eye cluster, the results of remaining landmarks from other clusters become worse slightly. This is because different facial parts have correlations and the occlusions of the left eye influences results of other facial parts. Note that WM and AM still perform well on occluded left eyes with the mean error of 6.60 and 6.50 respectively, due to the following reasons. First, WM weights each landmark proportional to its alignment error, which exploits correlations among landmarks. Second, AM uses an independent shape prediction layer focusing on a certain cluster of landmarks with small weights $u_j > 0$, $j \in Q^i$ in Eq. 9 for remaining landmarks, respectively, where correlations among landmarks are further exploited.

### E. Weighting Fine-Tuning for State-of-the-Art Frameworks

Most recently, there are a few well-designed and well-trained deep learning frameworks advancing the performance of face alignment, in which DAN [45] is a typical work. DAN uses cascaded deep neural networks to refine the localization accuracy of landmarks iteratively, where the entire face image and the landmark heatmap generated from the previous stage

TABLE VII
MEAN ERROR RESULTS OF WM AND AM ON COFW W/O OCCLU. AND W/ OCCLU. RESPECTIVELY. MEAN ERROR OF LANDMARKS FROM THE LEFT EYE CLUSTER, AND MEAN ERROR OF REMAINING LANDMARKS FROM OTHER CLUSTERS ARE BOTH SHOWN.

| Method | Left Eye | | Others | |
|---|---|---|---|---|
| | w/o occlu. | w/ occlu. | w/o occlu. | w/ occlu. |
| WM | 5.98 | 6.60 | 6.17 | 6.30 |
| AM | 5.88 | 6.50 | 6.06 | 6.18 |

are used in each stage. To evaluate the effectiveness of our method extended to state-of-the-art frameworks, we conduct experiments with our proposed weighting fine-tuning being applied to DAN. In particular, each stage of DAN is first pre-trained and further weighting fine-tuned, where DAN with weighting fine-tuning is named DAN-WM.

TABLE VIII
RESULTS OF MEAN ERROR OF DAN, RE-DAN, AND DAN-WM ON IBUG.

| Method | DAN [45] | re-DAN | **DAN-WM** |
|---|---|---|---|
| IBUG | 7.57 | 7.97 | 7.81 |

Note that the results of retrained DAN (re-DAN) using the published code[45] are slightly worse than reported results of DAN [45]. For a fair comparison, the results of mean error of DAN, re-DAN, and DAN-WM on IBUG benchmark are all shown in Table VIII. It can be seen that the mean error of re-DAN is reduced from 7.97 to 7.81 after using our proposed weighting fine-tuning. Note that our method uses only a single neural network, which has a concise structure with low model complexity. Our network can be replaced with a more powerful one such as cascaded deep neural networks, which could further improve the performance of face alignment.

## V. CONCLUSION

In this paper, we have developed a novel multi-center learning framework with multiple shape prediction layers for face alignment. The structure of multiple shape prediction layers is beneficial for reinforcing the learning process of each cluster of landmarks. In addition, we have proposed the model assembling method to integrate multiple shape prediction layers into one shape prediction layer so as to ensure a low model complexity. Extensive experiments have demonstrated the effectiveness of our method including handling complex occlusions and appearance variations. First, each component of our framework including Global Average Pooling, multiple shape prediction layers, weighting fine-tuning, and multi-center fine-tuning contributes to face alignment. Second, our proposed neural network and model assembling method allow real-time performance. Third, we have extended our method for detecting partially occluded faces and integrating with state-of-the-art frameworks, and have shown that our method exploits correlations among landmarks and can further improve the performance of state-of-the-art frameworks. The proposed framework is also promising to be applied for other face analysis tasks and multi-label problems.

## REFERENCES

[1] C. Cao, Q. Hou, and K. Zhou, "Displaced dynamic expression regression for real-time facial tracking and animation," *ACM Transactions on Graphics*, vol. 33, no. 4, p. 43, 2014.

[2] X. Fan, Z. Chai, Y. Feng, Y. Wang, S. Wang, and Z. Luo, "An efficient mesh-based face beautifier on mobile devices," *Neurocomputing*, vol. 172, pp. 134–142, 2016.

[3] B. Leng, Y. Liu, K. Yu, S. Xu, Z. Yuan, and J. Qin, "Cascade shallow cnn structure for face verification and identification," *Neurocomputing*, vol. 215, pp. 232–240, 2016.

[4] X. Cao, Y. Wei, F. Wen, and J. Sun, "Face alignment by explicit shape regression," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 2887–2894.

[5] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2013, pp. 532–539.

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

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

[8] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2013, pp. 3476–3483.

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

[10] 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.

[11] 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.

[12] ——, "Learning deep representation for face alignment with auxiliary attributes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 5, pp. 918–930, 2016.

[13] Z. Shao, H. Zhu, Y. Hao, M. Wang, and L. Ma, "Learning a multi-center convolutional network for unconstrained face alignment," in *IEEE International Conference on Multimedia and Expo*. IEEE, 2017, pp. 109–114.

[14] M. Lin, Q. Chen, and S. Yan, "Network in network," in *International Conference on Learning Representations*, 2014.

[15] 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.

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

[17] M. Pedersoli, T. Tuytelaars, and L. Van Gool, "Using a deformation field model for localizing faces and facial points under weak supervision," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2014, pp. 3694–3701.

[18] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.

[19] D. Lee, H. Park, and C. D. Yoo, "Face alignment using cascade gaussian process regression trees," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2015, pp. 4204–4212.

[20] P. Nair and A. Cavallaro, "3-d face detection, landmark localization, and registration using a point distribution model," *IEEE Transactions on Multimedia*, vol. 11, no. 4, pp. 611–623, 2009.

[21] 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 *IEEE International Conference on Computer Vision*. IEEE, 2013, pp. 1944–1951.

[22] 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.

[23] X. Yu, Z. Lin, J. Brandt, and D. N. Metaxas, "Consensus of regression for occlusion-robust facial feature localization," in *European Conference on Computer Vision*. Springer, 2014, pp. 105–118.

[24] Y. Wu and Q. Ji, "Robust facial landmark detection under significant head poses and occlusion," in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 3658–3666.

[25] S. Honari, J. Yosinski, P. Vincent, and C. Pal, "Recombinator networks: Learning coarse-to-fine feature aggregation," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 5743–5752.

[26] Z. Shao, S. Ding, Y. Zhao, Q. Zhang, and L. Ma, "Learning deep representation from coarse to fine for face alignment," in *IEEE International Conference on Multimedia and Expo*. IEEE, 2016, pp. 1–6.

[27] S. Xiao, J. Feng, J. Xing, H. Lai, S. Yan, and A. Kassim, "Robust facial landmark detection via recurrent attentive-refinement networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 57–72.

[28] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2015.

[29] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning*, 2015, pp. 448–456.

[30] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *International Conference on Machine Learning*, 2010, pp. 807–814.

[31] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.

[32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 770–778.

[33] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.

[34] M. Köstinger, P. Wohlhart, P. M. Roth, and H. Bischof, "Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization," in *IEEE International Conference on Computer Vision Workshops*. IEEE, 2011, pp. 2144–2151.

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

[36] S. Milborrow, J. Morkel, and F. Nicolls, "The muct landmarked face database," *Pattern Recognition Association of South Africa*, vol. 201, no. 0, 2010.

[37] P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar, "Localizing parts of faces using a consensus of exemplars," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 12, pp. 2930–2940, 2013.

[38] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," Technical Report 07-49, University of Massachusetts, Amherst, Tech. Rep., 2007.

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

[40] V. Le, J. Brandt, Z. Lin, L. Bourdev, and T. S. Huang, "Interactive facial feature localization," in *European Conference on Computer Vision*. Springer, 2012, pp. 679–692.

[41] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *ACM International Conference on Multimedia*. ACM, 2014, pp. 675–678.

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

[43] Z. Shao, S. Ding, H. Zhu, C. Wang, and L. Ma, "Face alignment by deep convolutional network with adaptive learning rate," in *IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2016, pp. 1283–1287.

[44] Y. Wu, C. Gou, and Q. Ji, "Simultaneous facial landmark detection, pose and deformation estimation under facial occlusion," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3471–3480.

[45] M. Kowalski, J. Naruniec, and T. Trzcinski, "Deep alignment network: A convolutional neural network for robust face alignment," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, vol. 3, 2017, p. 6.