

MobileFAN: Transferring Deep Hidden Representation for Face Alignment

Yang Zhao^{1,2} Yifan Liu² Chunhua Shen² Yongsheng Gao¹ Shengwu Xiong³

¹Griffith University ²The University of Adelaide ³Wuhan University of Technology

Abstract

Facial landmark detection is a crucial prerequisite for many face analysis applications. Deep learning-based methods currently dominate the approach of addressing the facial landmark detection. However, such works generally introduce a large number of parameters, resulting in high memory cost. In this paper, we aim for a lightweight as well as effective solution to facial landmark detection. To this end, we propose an effective lightweight model, namely Mobile Face Alignment Network (MobileFAN), using a simple backbone MobileNetV2 as the encoder and three deconvolutional layers as the decoder. The proposed MobileFAN, with only 8% of the model size and lower computational cost, achieves superior or equivalent performance compared with state-of-the-art models. Moreover, by transferring the geometric structural information of a face graph from a large complex model to our proposed MobileFAN through feature-aligned distillation and feature-similarity distillation, the performance of MobileFAN is further improved in effectiveness and efficiency for face alignment. Extensive experiment results on three challenging facial landmark estimation benchmarks including COFW, 300W and WFLW show the superiority of our proposed MobileFAN against state-of-the-art methods.

1. Introduction

Facial landmark detection, a.k.a, face alignment, is a crucial step for various downstream face applications including face recognition [27], facial attributes estimation [45], face pose estimation [14] and so forth. Face alignment aims to find the coordinates of several predefined landmarks or parts, such as eye center, eyebrow, nose tip, mouth and chin, on a face graph. Although great progress has been made on accuracy improvements in the past decades [43, 15], approaches focusing on simple, small and lightweight networks for face alignment receive relatively much less attention.

Significant improvements via deep Convolutional Neural Networks (CNNs) have been achieved on facial landmark detection recently [30, 4, 29], even though it remains a very challenging task when dealing with faces in

real-world conditions (*e.g.*, faces with unconstrained large pose variations and heavy occlusions). In order to guarantee promising performance in face alignment benchmarks, the majority of those works are designed to adopt large backbones (*e.g.*, Hourglass [21] and ResNet-50 [10]), carefully designed schemes (*e.g.*, a coarse-to-fine cascade regression framework [28]), or adding extra face structure information (*e.g.*, face boundary information [30]). Recently, neural networks with small model size, light computation cost and high accuracy, have attracted much attention because of the need of applications on mobile devices [38]. In this paper, we aim to investigate the possibility of optimizing facial landmark detection with a simpler and smaller model. We propose a plain model without bells and whistles, namely *Mobile Face Alignment Network* (MobileFAN), which employs an Encoder-Decoder architecture in the form of Convolution-Deconvolution Network. In the proposed MobileFAN, MobileNetV2 [26] is adopted as the encoder, while the decoder is constructed utilizing three deconvolutional layers. Model details are illustrated in Section III.

More recently, knowledge distillation (KD) has attracted much attention for its simplicity and efficiency [11]. Motivated by KD and TCNN [33], which has shown that intermediate features from deep networks are good at predicting different head poses in facial landmark detection, we further introduce knowledge transfer techniques to help the training of our proposed lightweight face alignment network. Because our proposed MobileFAN uses several deconvolutional layers sequentially to map from the input space to the output space, we try to transfer the useful information of intermediate feature maps from a teacher to a student.

Inspired by [18, 24], we propose to align the deconvolutional feature maps between student models and teacher models. Specifically, the feature map generated by the student network can be transformed to a new feature map, which needs to match the same size of the corresponding feature map generated by the teacher network. Mean squared error (MSE) is used as the loss function to measure the distance between teacher's feature map and student's new feature map. We term this scheme feature-aligned distillation, which can transfer the distribution of intermediate

feature map produced by the teacher network to that of the student network.

To distill more structured knowledge information from the teacher network, inspired by [19], we apply the feature-similarity distillation to our framework. The similarity matrix is generated by computing cosine similarity of feature vectors. We find that the similarity matrix can be used to represent the structure information of a face image. It contains the directional knowledge between features, which can be thought of a kind of structure information. With the help of feature-similarity distillation, the student network is trained to make its similarity matrix similar to that of the teacher network. The illustration of our method for knowledge transfer is depicted in Fig. 1 (c).

The interest of this work lies in exploring a simple, small and lightweight network that can achieve comparable or even better results than the common facial landmark detection benchmarks. To summarize, our main contributions are as follows.

- We propose a simple and lightweight *Mobile Face Alignment Network* (MobileFAN) for face alignment. The proposed MobileFAN achieves comparable results, while having less than 8% of parameters compared with the sizes of state-of-the-arts.
- Feature-aligned distillation and feature-similarity distillation are introduced and integrated with the proposed MobileFAN, leading to further improvements in alignment accuracy.
- Extensive experimental results on three challenging benchmark datasets demonstrate the efficiency and effectiveness of our method over state-of-the-art methods in face alignment.

2. Related Work

In this section, we present an overview of related work on facial landmark detection and knowledge distillation.

2.1. Facial Landmark Detection

Traditional Methods. Facial landmark detection has been an active topic for more than twenty years. Recently, cascade regression attracts a lot of attention, which focuses on learning a cascade of regressors to iteratively update the shape estimation. BurgosArtizzu *et al.* [1] proposed Robust Cascade Pose Regression (RCPR) that reduces exposure to outliers by detecting occlusions explicitly and using robust shape-indexed features. Explicit Shape Regression (ESR) [2] introduced two-level boosted regression and correlation-based feature selection. Ren *et al.* [22] proposed Local Binary Features (LBF) that is computationally cheap and thus enables very fast regression on the face alignment tasks.

CNN-based Methods. Apart from the above early works, deep learning-based face alignment approaches have achieved state-of-the-art performance. They can be divided into two categories, namely coordinate regression-based method and heatmap regression-based method. A coordinate regression-based method estimates the landmark coordinates vector from the input image directly. The earliest work could be dated to [28]. Sun *et al.* [28] trained a three-level cascade CNN to locate the facial landmarks in a coarse-to-fine manner, and obtained promising landmark detection results. A multi-task learning framework is proposed by Zhang *et al.* [42] to optimize face alignment and correlated facial attributes, such as pose, expression and gender, simultaneously. More recently, Feng *et al.* [7] proposed a new loss function, namely Wingloss, to fill the gap of a better loss function in facial coordinates regression community. It shows that Wingloss with the proposed strong data augmentation method, pose-based data balancing (PDB), could obtain better performance against widely used L2 loss. Different from the above methods, our approach regards face alignment as a dense prediction problem.

A heatmap regression-based method generates a probability heatmap for each landmark respectively. Thanks to the development of Hourglass [21], heatmap regression has been successfully applied to landmark localization problems. Yang *et al.* [37] adopted a supervised face transformation based on the Hourglass to reduce the variance of the target. LAB [30] utilized boundary lines to characterize the geometric structure of a face image and thus improved the detection of facial landmarks. However, both the two methods rely on the Hourglass, resulting in introducing a large number of parameters. Valle *et al.* [29] used a simple CNN to generate heatmaps of landmark locations for a better initialization to Ensemble of Regression Trees (ERT) regressor.

By contrast, our model requires neither cascaded networks nor large backbones, leading to great reduction in model parameters and computation complexity, whilst still achieving comparable or even better accuracy.

2.2. Knowledge Distillation

Deep CNN models dominate the approach to solving many computer vision tasks recently [13, 39]. However, millions of parameters are commonly introduced in these Deep CNN models, leading to large model sizes and expensive computation cost. As a result, it is difficult to deploy such models to real-time applications. Therefore, it motivates researchers to focus on smaller networks that can fit large training data while maintaining the performance. Recently, knowledge distillation (KD) [11] has attracted much attention due to its capability of transferring rich information from a large and complex teacher network to a small

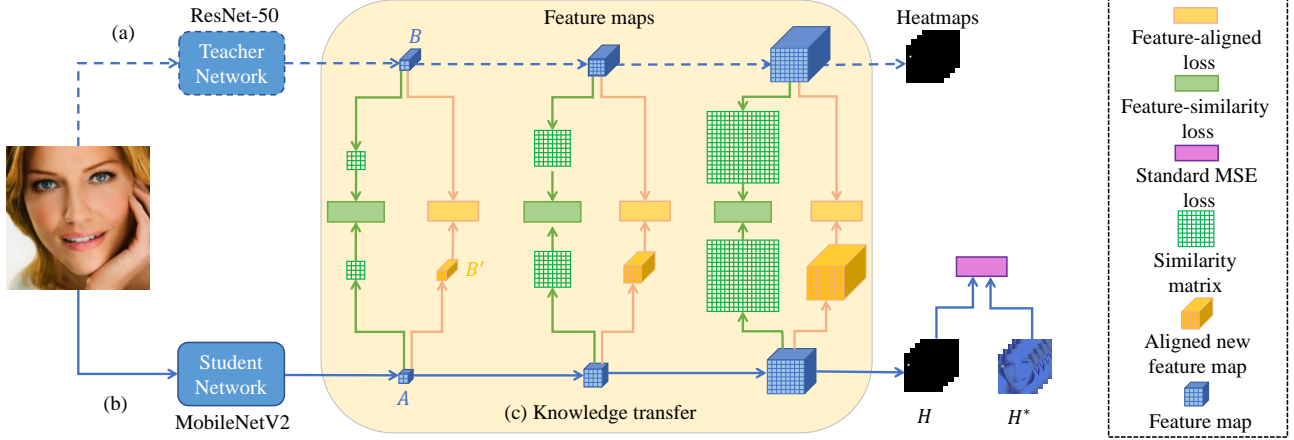


Figure 1: The structure of proposed MobileFAN and an overview of our knowledge transfer framework. (a) The process indicated by the blue dotted arrow is the teacher network, which is made up of ResNet-50 and three deconvolutional layers. (b) The process indicated by the blue arrow is MobileFAN (the student network), which consists of MobileNetV2 and three deconvolutional layers. (c) Knowledge transfer module, which is designed to help the training process of MobileFAN by introducing feature-aligned distillation and feature-similarity distillation.

and compact student network. It is widely used in model compression. Originally, KD is used in the task of image classification [11], where a compact model can learn from the output of a large model, namely soft target. So the student is supervised by both softened labels and hard labels simultaneously.

Following [11], some subsequent works have tried to transfer intermediate representations of the teacher network to that of the student network, and achieved great progress in image classification [24, 40], object detection [18], and semantic segmentation [19]. Romero *et al.* proposed a FitNet [24] which directly aligns full feature maps of the teacher model and student model. Attention transfer (AT) [40] was proposed to regularize the learning of the student network by imitating the attention maps of a powerful teacher network. Liu *et al.* [19] proposed to distill the pair-wise information from the teacher model to a student model via the last convolution features. Unlike previous approaches, we perform the distillation through multiple features.

Given the effectiveness of knowledge distillation in the above applications, we are motivated to simultaneously use both feature-aligned distillation and feature-similarity distillation on intermediate features in this work.

3. Method

In this section, we start with an introduction of network architectures. Then we take a look at standard MSE loss and two knowledge distillation schemes: feature-aligned distillation and feature-similarity distillation.

3.1. Network Architectures

Teacher Network. Following [34], we introduce a teacher network which employs an Encoder-Decoder architecture

in the form of Convolution-Deconvolution Network. The encoder is a ResNet50 where the last average pooling layer and the classification layer are both removed. The decoder consists of three deconvolutional layers, with the dimension of 256 and 4×4 kernel for each layer. The stride of each deconvolutional layer in the decoder is set to be 2. Then, a 1×1 convolutional layer is added after the decoder to generate likelihood heatmaps. The architecture of the teacher network is shown in Fig. 1 (a).

Student Network. To make our framework easy to be reconstructed, we present a student network using a similar Encoder-Decoder structure as adopted in the teacher network. The encoder is a MobileNetV2 [26] with the last three layers (*i.e.*, one global average pooling layer and two convolutional layers) being removed. The decoder is composed of three deconvolutional layers, and is added over the last bottleneck of the MobileNetV2. Each deconvolutional layer has 128 filters with 2×2 kernel. The stride of each deconvolutional layer in the student network is 2. Same as in [21, 34], a 1×1 convolutional layer is added after the decoder to generate likelihood heatmaps $\mathbf{H} = \{\mathbf{H}_l\}_{l=1}^L$ for L facial landmarks. An illustration of the detailed student network structure is shown in Fig. 1 (b).

3.2. Loss Function

Mean Squared Error. Same as [34], we employ Mean Squared Error (MSE) loss to compare the predicted heatmaps \mathbf{H} and the ground-truth heatmaps \mathbf{H}^* generated from the annotated 2D facial landmarks.

Specifically, $\mathbf{H}^* = \{\mathbf{H}_l^*\}$ is a set of L response maps, one per facial landmark, where $\mathbf{H}_l^* \in \mathbb{R}^{64 \times 64}$, $l \in \{1 \dots L\}$. Here heatmap \mathbf{H}_l^* for the l^{th} landmark is made up of a 2D gaussian centered on the landmark location. Let $x \in \mathbb{R}^2$ be the ground-truth position of the l^{th} facial landmark, the

value at location $p \in \mathbb{R}^2$ in \mathbf{H}_l^* is defined as:

$$\mathbf{H}_l^*(p) = \exp\left(-\frac{\|p - x\|_2^2}{2\sigma^2}\right). \quad (1)$$

Therefore, the loss between the predicted heatmaps \mathbf{H} and ground-truth heatmaps \mathbf{H}^* is defined as:

$$\mathcal{L}(\mathbf{H}, \mathbf{H}^*) = \|\mathbf{H} - \mathbf{H}^*\|^2. \quad (2)$$

After training, the l^{th} landmark location can be generated from the corresponding predicted heatmap \mathbf{H}_l by transforming the highest heatvalued location from $1/4$ to the original image space.

Knowledge Transfer. In the student-teacher framework, apart from the standard MSE loss, $\mathcal{L}(\mathbf{H}, \mathbf{H}^*)$, we further introduce knowledge transfer loss to help the training of our student network (MobileFAN). In other words, we want the student network to learn not only the information provided by the ground-truth labels, but also the finer structure knowledge encoded by the teacher network. Let \mathcal{T} , \mathcal{S} and $\mathbf{W}_{\mathcal{T}}$, $\mathbf{W}_{\mathcal{S}}$ denote the teacher network, the student network and their corresponding weights. Details of knowledge distillation are described below.

Feature-Aligned Distillation. In order to transfer richer facial details (*e.g.*, exaggerated expressions and head poses) learned by a teacher network to the student network, we perform feature-aligned distillation such that the distribution of a feature of the student network is similar to that of the teacher network. Feature-aligned distillation is designed to align the feature map between a student network and a teacher network. Given a feature map $\mathbf{A} \in \mathbb{R}^{C \times H \times W}$ from the student network \mathcal{S} and a feature map $\mathbf{B} \in \mathbb{R}^{C' \times H \times W}$ from the teacher network \mathcal{T} , where C and C' represent channels and $H \times W$ is the spatial dimensions. To perform the feature-aligned distillation, \mathbf{A} and \mathbf{B} should have the same size (*i.e.*, $C = C'$ in our case). Therefore, we firstly adopt a 1×1 convolutional layer to take an input of feature map \mathbf{A} and output a new feature map $\mathbf{B}' \in \mathbb{R}^{C' \times H \times W}$ that has the same size of the feature map \mathbf{B} in the teacher network. Therefore, the feature-aligned transfer loss between the teacher network and student network is defined as:

$$\mathcal{L}_{FA} = \|\mathbf{B} - \mathbf{B}'\|^2. \quad (3)$$

Feature-Similarity Distillation. Facial images are geometrically constrained. As illustrated before, we adopt feature-similarity distillation to transfer more structural information from the teacher network to the student network (MobileFAN) by comparing their similarity matrix. The similarity matrix represents the basic facial structures and textures, which can provide richer directional information to facial landmark detection. We perform cosine similarity computation on the whole feature map, making the relative spatial positions between facial landmarks more precisely.

Given a feature map that has dimensions of $C \times H \times W$, where C is the total number of channels and $H \times W$ is the feature map size, $\mathbf{f}_i \in \mathbb{R}^C$ denotes a feature vector extracted from the i^{th} ($i = 1, \dots, H \times W$) spatial location of this feature map. Therefore, the cosine similarity a_{ij} between the i^{th} feature vector \mathbf{f}_i and j^{th} feature vector \mathbf{f}_j is calculated as:

$$a_{ij} = \frac{\mathbf{f}_i \cdot \mathbf{f}_j}{\|\mathbf{f}_i\| \|\mathbf{f}_j\|}. \quad (4)$$

Let a_{ij}^s denote the similarity between the i^{th} feature vector and j^{th} feature vector computed from the feature map \mathbf{A} , and a_{ij}^t denote the similarity between the i^{th} feature vector and j^{th} feature vector computed from the feature map \mathbf{B} . The feature-similarity transfer loss is then formulated as:

$$\mathcal{L}_{FS} = \frac{\sum_{i \in K} \sum_{j \in K} (a_{ij}^s - a_{ij}^t)^2}{(H \times W)^2}, \quad (5)$$

where $K = \{1, 2, \dots, H \times W\}$ denotes all the locations.

3.3. Distillation over Scales

In order to transfer low-, mid- and high-level useful geometric structure information from the teacher network to the student network (MobileFAN), we extend the combination of the feature-aligned distillation and feature-similarity distillation to three deconvolutional layers. As shown in Fig. 1 (c), three deconvolutional layers of the student network are guided by that of the teacher network, where significantly richer facial details are provided. It is different from most previous methods [19] that only add supervision on the last convolutional layer. Then we perform both feature-aligned distillation and feature-similarity distillation on the three deconvolutional feature maps during training. So the student network is trained to optimize the following loss function:

$$\mathcal{L}_{KD} = \mathcal{L}(\mathbf{H}, \mathbf{H}^*) + \lambda \sum_{r=1}^3 (\mathcal{L}_{FA}^r + \mathcal{L}_{FS}^r), \quad (6)$$

where λ is a tunable parameter to balance the MSE loss and the distillation loss. \mathcal{L}_{FA}^r and \mathcal{L}_{FS}^r are the feature-aligned loss and feature-similarity loss of the r^{th} deconvolutional layer. Extensive experiments show that with the help of distilled knowledge from different feature map scales, the performance of facial landmark detection can be significantly increased.

3.4. Learning Procedure

Training the proposed MobileFAN. To evaluate the performance of the proposed *vanilla* MobileFAN, we optimize MobileFAN only with the standard MSE loss (as illustrated

in Equation (2)) without any extra losses. The experimental results indicate that our proposed MobileFAN, a simple and small network, still can handle the problem of facial landmark detection with satisfying performance.

Training MobileFAN with distilled knowledge. To transfer the distilled knowledge from a large complicated network to the proposed MobileFAN, we regard MobileFAN as a student network in a student-teacher framework. Fig. 1 summarizes the training of the knowledge transfer framework. Specifically, a teacher network (Fig. 1 (a)) is pre-trained and the parameters are kept frozen during training. The training stage of the proposed MobileFAN is supervised by standard MSE loss, feature-aligned loss and feature-similarity loss. In other words, guided by the pre-trained parameters \mathbf{W}_T of the teacher network, we train the parameters of the MobileFAN \mathbf{W}_S to minimize Equation (6).

4. Experiments

4.1. Datasets

We perform experiments on three challenging public datasets: the Caltech Occluded Faces in the Wild (COFW) dataset [1], the 300 Faces in the Wild (300W) dataset [25] and the Wider Facial Landmarks in the Wild (WFLW) dataset [30].

COFW. The face images in COFW comprise heavy occlusions and large shape variations, which are common issues in realistic conditions [1]. Its training set has 1345 faces, and the testing set has 507 faces. Each image in the COFW dataset has 29 manually annotated landmarks, as shown in Fig. 2(a).

300W. The 300W [25] dataset is a widely used facial landmark detection benchmark, which consists of HELEN, LFPW, AFW and IBUG datasets. Images in HELEN, LFPW and AFW datasets are collected in the wild environment, where large pose variations, expression variations, and partial occlusions may exist. There are 68 annotated facial landmarks in each face from 300W dataset, as shown in Fig. 2(b). We follow the same protocol as used in [23] to adopt 3148 images for training (2000 images from the training subset of HELEN dataset, 811 images from the training subset of LFPW dataset and 337 images from the full set of AFW dataset). For testing, the Full test set has 689 images including Common subset (554 images) and Challenging subset (135 images). Here Common subset is composed of HELEN test subset (330 images) and LFPW test subset (224 images), while Challenging subset is the IBUG dataset.

WFLW. WFLW [30] is a recently proposed facial landmark dataset based on WIDER FACE. It comprises 7500 face images (for training) and 2500 face images (for testing) with 98 manual annotated landmarks (shown in Fig. 2(c)), respectively. Faces in WFLW are collected under uncon-

Table 1: A summary of the evaluation protocols used in our experiments.

Protocol	Training Set Size	Test Set Size	#Landmarks	Normalisation Term
COFW	1, 345	507	29	inter-ocular distance
300W Full set	3, 148	689	68	inter-ocular distance
WFLW	7, 500	2, 500	98	inter-ocular distance

strained conditions, such as large variations in poses, exaggerated expressions and heavy occlusions. To validate the robustness against each different condition, WFLW is further divided into several subsets including large pose (326 images), expression (314 images), illumination (698 images), make-up (206 images), occlusion (736 images) and blur (773 images). We report the results of all the competing methods on the whole test set and each testing subset in the WFLW dataset.

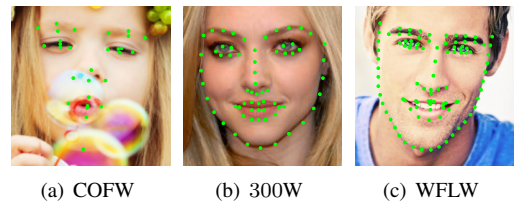


Figure 2: An illustration of the landmark annotations of (a) COFW dataset, (b) 300W dataset and (c) WFLW dataset.

4.2. Evaluation Metrics and Implementation Details

Evaluation Metrics. We adopt the normalized mean error and the area-under-the curve (AUC) as metrics for evaluation. For all the datasets (COFW, 300W and WFLW), we use the distance between the outer eye corners (inter-ocular distance) as the normalization term [1, 30].

The mean error is defined as the average Euclidean distance between the predicted facial landmark locations $p_{i,j}$ and their corresponding ground-truth facial landmark locations $g_{i,j}$:

$$error = \frac{1}{N} \sum_{i=1}^L \frac{\frac{1}{L} \sum_{j=1}^L |p_{i,j} - g_{i,j}|^2}{d}. \quad (7)$$

where N is the number of images in the test set, and L is the number of landmarks (as illustrated in TABLE 1). d is the normalization factor, which is the inter-ocular distance.

We also measure the Cumulative Errors Distribution (CED) curve and the failure rate (which is defined as the proportion of failed detected faces) on these benchmarks. Specifically, any normalized error above 0.1 is considered as a failure [30]. The summary of detailed evaluation protocols used in our experiments is listed in TABLE 1.

Implementation Details. All the face images including both training and testing images are cropped and scaled to 256×256 according to center location and provided bounding box [30, 3]. Standard data augmentation is performed to make networks robust to data variations. Specifically,

Table 2: Mean error (%) and failure rate (%) on COFW test set.

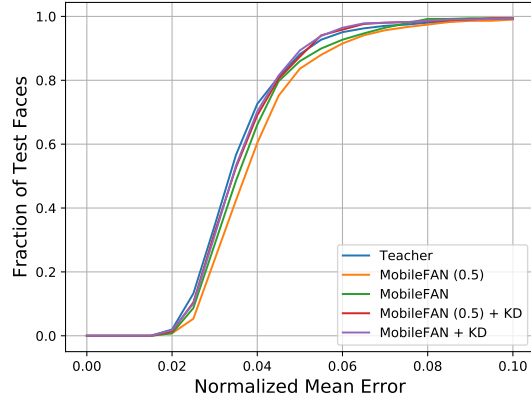
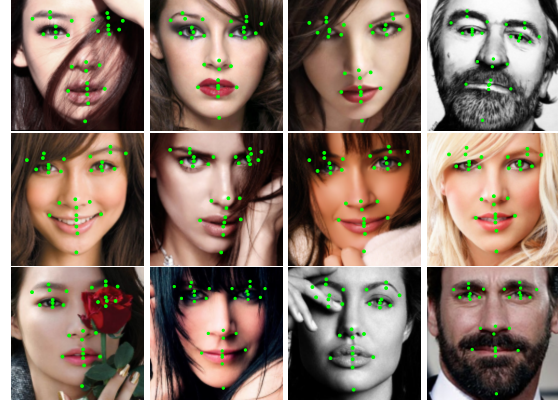
Method	Mean Error	Failure Rate
Human [1]	5.6	-
RCPR [1]	8.50	20.00
HPM [9]	7.50	13.00
CCR [6]	7.03	10.9
DRDA [41]	6.46	6.00
RAR [35]	6.03	4.14
SFPD [32]	6.40	-
DAC-CSR [8]	6.03	4.73
CNN6 (Wing + PDB) [7]	5.44	3.75
ResNet50 (Wing + PDB) [7]	5.07	3.16
LAB [30]	3.92	0.39
Teacher	3.65	0.59
MobileFAN (0.5)	4.01	0.99
MobileFAN	3.82	0.59
MobileFAN (0.5) + KD	3.68	0.59
MobileFAN + KD	3.66	0.59

we follow [7, 30] to augment samples by (± 30 degree) in-plane rotation, (0.75-1.25) scaling and randomly horizontal flip with the probability of 50%. In the training, Adam optimizer is used with a mini-batch size of 8 for 80 epochs. The base learning rate is 10^{-3} , and it drops to 10^{-4} at 30th epochs and 10^{-5} at 50th epochs respectively. λ is set to be 10^{-2} for COFW dataset and 10^{-4} for 300W and WFLW dataset. For implementation, all our face estimation models are trained with Pytorch 0.4.0 toolbox on one GPU.

4.3. Comparison with state-of-the-art methods

We compare the proposed method against the state-of-the-art methods on each dataset. To further explore whether the effectiveness of a smaller decoder on face alignment task, we adopt a channel-halved version of MobileFAN, where we use 64 dimension to replace 128 dimension of each deconvolutional layer. This architecture is denoted as MobileFAN (0.5). We apply our distillation method to both of the two lightweight networks: MobileFAN and MobileFAN (0.5). For simplicity, we name our full models trained using the combination of feature-aligned distillation and feature-similarity distillation of all the deconvolutional layers to be “MobileFAN + KD” and “MobileFAN (0.5) + KD”. Similarly, the baseline models, MobileFAN without distillation and MobileFAN (0.5) without distillation, are named as “MobileFAN” and “MobileFAN (0.5)”. We use “Teacher” to represent our proposed teacher network.

Evaluation on COFW. In TABLE 2, we provide the results of the state-of-the-art methods in COFW test set. We can see that the proposed simple and small “MobileFAN” achieves 3.82% mean error with 0.59% failure rate without any extra information. Although the mean error of “MobileFAN (0.5)” is a little higher than that of LAB [30], it is still a comparable result (4.01% mean error with 0.99% failure rate). With the knowledge transferred from “Teacher”, our method outperforms existing methods with a margin of over 0.24% in mean error reduction. More specifically, our proposed models, “MobileFAN + KD” and “MobileFAN (0.5) + KD”, achieve the best performances on COFW dataset, with about 6.63% and 6.12% improvements in mean error

**Figure 3:** CED curves of the baselines, teacher network and the proposed MobileFAN with KD on COFW test set.**Figure 4:** Example alignment results of MobileFAN on COFW [1] test set.

reduction over LAB [30] with extra boundary information. “MobileFAN + KD” achieves the comparable result (3.66% mean error with 0.59% failure rate) to the teacher network (3.65% mean error with 0.59% failure rate). This is not surprising since our proposed distillation method provides rich structural information of a face image, which may contribute to the performance of facial landmark detection. The CED curves in Fig. 3 show that the distilled small networks gain better performance than baselines as well as achieve comparable performance to the “Teacher” network. Fig. 4 shows some example alignment results, demonstrating the effectiveness on various occlusions of “MobileFAN”.

Evaluation on 300W. The 300W dataset is a challenging face alignment benchmark because of its variants on pose and expressions. TABLE 3 shows the comparable performance with previous methods on 300W dataset. We can observe that simple “MobileFAN” performs better than the state-of-the-art SAN [4], but the number of model parameters of “MobileFAN” is $28\times$ smaller than that of SAN (we can see from TABLE 5). Although “MobileFAN + KD” does not outperform DCFE [29], it achieves comparable re-

Table 3: Mean error (%) on 300W Common subset, Challenging subset and Full set.

Method	Common	Challenging	Full
RCN [12]	4.67	8.44	5.41
DAN [16]	3.19	5.24	3.59
PCD-CNN [17]	3.67	7.62	4.44
CPM [5]	3.39	8.14	4.36
DSRN [20]	4.12	9.68	5.21
SAN [4]	3.34	6.60	3.98
LAB [30]	2.98	5.19	3.49
DCFE [29]	2.76	5.22	3.24
Teacher	2.97	5.23	3.41
MobileFAN (0.5)	5.44	8.24	5.99
MobileFAN	3.10	5.62	3.60
MobileFAN (0.5) + KD	4.22	6.87	4.74
MobileFAN + KD	2.98	5.34	3.45

sults to LAB [30] with extra boundary information on 300W Full set and Common subset. Using the knowledge distillation, our two full models are better than their corresponding baselines. The “MobileFAN + KD” achieves 4.17%, 4.98% and 3.81% improvements over “MobileFAN” on 300W Full set, Challenging subset and Common subset. Although the “MobileFAN (0.5) + KD” fails to compete other state-of-the-art methods, which is possible because the dimension of output score maps (for 300W dataset, it is 68) is larger than the dimension of the final deconvolutional layer (64), it reduces the mean error from 5.99% to 4.74% on 300W Full set over its baseline “MobileFAN (0.5)”. Fig. 5 visualizes some of our results. It can be observed that driven by knowledge transfer technique, our model can capture various facial expressions accurately.

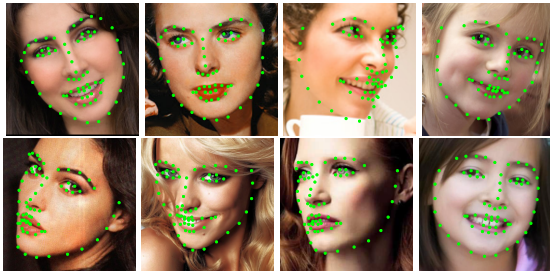


Figure 5: Example alignment results of MobileFAN + KD on 300W [25] Full set.

Evaluation on WFLW. A summary of the performance obtained by state-of-the-art methods and the proposed approach on WFLW test set and six subset is shown in TABLE 4. As indicated in TABLE 4, our proposed “MobileFAN” outperforms LAB [30] with boundary information in test set and all six subset and ResNet50 (Wing + PDB) [7] with strong data augmentation in Test set, Make-up and Occlusion subset. Although “MobileFAN” performs a little worse than ResNet50 (Wing + PDB) [7] in remaining subset, it achieves comparable results with merely 8% of parameters of ResNet50 (Wing + PDB) (which can be observed from TABLE 5). We can see that “MobileFAN + KD” model outperforms the-state-of-art methods, with a mean error of 4.93% and failure rate of 5.32% on WFLW test set. In particular, compared with former best models, ResNet50 (Wing + PDB) [7] and LAB [30], our “Mobile-

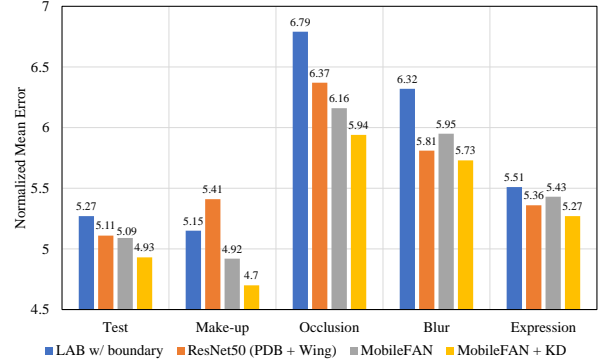


Figure 6: Normalized Mean Error (%) on WFLW [30] Test set and 4 typical subset for the method of LAB [30], ResNet50 (Wing + PDB) [7] and the proposed method.

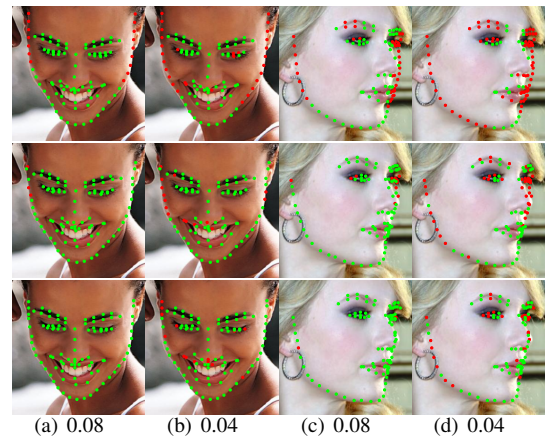


Figure 7: Results of example samples of WFLW [30] under different mean error thresholds: 0.08 and 0.04. Top row: Results of LAB [30]. Middle row: Results of MobileFAN. Bottom row: Results of MobileFAN + KD. Red points indicate normalized mean error is larger than threshold.

FAN + KD” achieves significant mean error reduction with respect to ResNet50 (Wing + PDB) [7] and LAB [30] of 3.52% and 6.45% on WFLW test set, respectively. Similarly, the failure rate is reduced from 6.00% to 5.32% and from 7.56% to 5.32% compared with ResNet50 (Wing + PDB) [7] and LAB [30].

To provide a more straightforward comparative illustration, we compare in Fig. 6 our “MobileFAN” and “MobileFAN + KD” against ResNet50 (Wing + PDB) and LAB on WFLW test set and four typical subsets. We can see that “MobileFAN” outperforms ResNet50 (Wing + PDB) and LAB with a big margin on WFLW Test set, Make-up subset and Occlusion subset, let along “MobileFAN + KD” with the help of the “Teacher”. In particular, “MobileFAN” achieves 9.06% relative improvement in mean error reduction over ResNet50 (Wing + PDB) on Make-up subset, as well as 9.28% relative improvement in mean error reduction over LAB on Occlusion subset. Although “MobileFAN” achieves comparable results compared with ResNet50 (Wing + PDB) on Blur subset and Expression subset, it outperforms LAB with a big margin by using less

Table 4: Mean error(%), failure rate (%) and AUC on WFLW test set and six subsets: pose, expression (expr.), illumination (illu.), make-up (mu.), occlusion (occu.) and blur.

Method	test	pose	expr.	illu.	mu.	occu.	blur
Mean Error							
ESR [2]	11.13	25.88	11.47	10.49	11.05	13.75	12.20
SDM [36]	10.29	24.10	11.45	9.32	9.38	13.03	11.28
CFSS [44]	9.07	21.36	10.09	8.30	8.74	11.76	9.96
DVLN [31]	6.08	11.54	6.78	5.73	5.98	7.33	6.88
LAB [30]	5.27	10.24	5.51	5.23	5.15	6.79	6.32
ResNet50(Wing+PDB) [7]	5.11	8.75	5.36	4.93	5.41	6.37	5.81
Teacher	4.82	8.48	5.07	4.77	4.55	5.88	5.59
MobileFAN (0.5)	5.70	10.17	5.95	5.68	5.47	6.66	6.49
MobileFAN (0.5) + KD	5.59	9.68	5.98	5.45	5.33	6.49	6.31
MobileFAN	5.09	8.97	5.43	5.03	4.92	6.16	5.95
MobileFAN + KD	4.93	8.72	5.27	4.93	4.70	5.94	5.73
Failure Rate							
ESR [2]	35.24	90.18	42.04	30.80	38.84	47.28	41.40
SDM [36]	29.40	84.36	33.44	26.22	27.67	41.85	35.32
CFSS [44]	20.56	66.26	23.25	17.34	21.84	32.88	23.67
DVLN [31]	10.84	46.93	11.15	7.31	11.65	16.30	13.71
LAB [30]	7.56	28.83	6.37	6.73	7.77	13.72	10.74
ResNet50(Wing+PDB) [7]	6.00	22.70	4.78	4.30	7.77	12.50	7.76
Teacher	4.28	20.25	2.55	3.58	3.40	8.56	6.21
MobileFAN (0.5)	7.04	35.89	6.69	6.45	7.77	12.36	7.89
MobileFAN (0.5) + KD	6.72	30.67	6.05	5.73	8.74	12.5	8.54
MobileFAN	5.8	26.99	4.46	5.01	6.80	11.01	8.67
MobileFAN + KD	5.32	23.93	4.46	5.01	6.80	10.33	7.24
AUC							
ESR [2]	0.2774	0.0177	0.1981	0.2953	0.2485	0.1946	0.2204
SDM [36]	0.3002	0.0226	0.2293	0.3237	0.3125	0.2060	0.2398
CFSS [44]	0.3659	0.0632	0.3157	0.3854	0.3691	0.2688	0.3037
DVLN [31]	0.4551	0.1474	0.3889	0.4743	0.4494	0.3794	0.3973
LAB [30]	0.5323	0.2345	0.4951	0.5433	0.5394	0.4490	0.4630
ResNet50(Wing+PDB) [7]	0.5504	0.3100	0.4959	0.5408	0.5582	0.4885	0.4918
Teacher	0.5414	0.2706	0.5215	0.5513	0.5534	0.4653	0.4775
MobileFAN (0.5)	0.4630	0.1605	0.4380	0.4758	0.4777	0.4029	0.4065
MobileFAN (0.5) + KD	0.4682	0.1728	0.4373	0.4819	0.4823	0.4028	0.4091
MobileFAN	0.5163	0.2270	0.4888	0.5276	0.5251	0.4417	0.4484
MobileFAN + KD	0.5296	0.2443	0.5039	0.5388	0.5442	0.4576	0.4661



Figure 8: Example face alignment results on WFLW [30] test set. Top row: Face alignment results of LAB [30]. Second row: Face alignment results of our proposed MobileFAN. Third row: Face alignment results of our MobileFAN + KD. Bottom row: Ground-truth annotations. Red points indicate normalized error is larger than 0.1.

number of parameters. With the help of knowledge distillation, “MobileFAN + KD” achieves the state-of-the-art performance on WFLW test set and six subsets. The results indicate that our proposed lightweight model is robust to extreme conditions. We can visually see the advantages of

our “MobileFAN” from Fig. 7. Specifically, we compare LAB, “MobileFAN” and “MobileFAN + KD” under different mean error thresholds. It can be found that the number of landmarks of low mean error of our method is more than that of LAB. And the third row in Fig. 7 depicts further im-

provements led by adding feature-aligned distillation and feature-similarity distillation, where the knowledge transfer techniques provide richer facial details to make relative spatial positions between facial landmarks more precisely.

Some example results of LAB¹, “MobileFAN”, “MobileFAN +KD” and ground-truth on WFLW test set is showed in Fig. 8. We can observe that “MobileFAN + KD” improves the accuracy of landmarks above the face contour (chin), eyebrow, eye corner and so on.

Model Size and Computational Cost Analysis. To further evaluate the model size and the computational complexity, we calculate the number of network parameters (#Params), the sum of floating point operations (FLOPs) and the speed of our approach and other competing methods. The FLOPs of our model is calculated on the resolution of 256×256 . Frames per second (fps) is adopted for measuring the computation speed. Here the fps calculation is performed on an NVIDIA GeForce GTX 1070 card. We notice that the model size of our compact network is the smallest. We can see from TABLE 5 that the proposed models have minimal parameters and lowest computation complexity against LAB [30] and ResNet50 (Wing + PDB) [7], while remaining effective for facial landmark localization. Specifically, MobileFAN (0.5) and MobileFAN just have 1.84M and 2.02M parameters, respectively. Although our MobileFAN has fewer parameters, *e.g.*, 8% of model size of LAB and ResNet50 (Wing + PDB) and 3.54% of model size of SAN, it achieves comparable or even better results competing with the state-of-the-art methods. It is observed that the proposed MobileFAN can process 238 fps, which is significantly faster than the 6, 50, and 126 fps of the state-of-the-art approaches.

Table 5: A comparison of different networks in backbone, model size (the number of model parameters), computational cost (FLOPS) and speed (fps).

Method	backbone	#Params (M)	FLOPs (B)	Speed(fps)
DVLN [31]	VGG-16	132.0	14.4	-
SAN [4]	ResNet-152	57.4	10.7	50
LAB [30]	Hourglass	25.1	19.1	6
ResNet50 (Wing + PDB) [7]	ResNet-50	25	3.8	126
MobileFAN	MobileNetV2	2.02	0.72	238
MobileFAN (0.5)	MobileNetV2	1.84	0.45	249

4.4. Ablation study

Our framework consists of several different components, such as feature-aligned distillation and feature-similarity distillation of different deconvolutional layers. In this section, we take a look at the effectiveness of different distillation methods on 300W Challenging subset. Based on the baseline network, MobileFAN without distillation (“MobileFAN”), we evaluate the mean error using various combinations of each component, as summarized in TABLE 6. In addition, we analyze the influence of the hyperparameter λ (as described in Section III) on COFW Test set.

Table 6: The proposed distillation components in our method.

Proposed distillation component	Abbreviation
feature-aligned distillation of the first deconvolutional layer	FA ₁
feature-aligned distillation of the second deconvolutional layer	FA ₂
feature-aligned distillation of the third deconvolutional layer	FA ₃
feature-similarity distillation of the first deconvolutional layer	FS ₁
feature-similarity distillation of the second deconvolutional layer	FS ₂
feature-similarity distillation of the third deconvolutional layer	FS ₃

Table 7: Mean error (%) of feature-aligned distillation of different deconvolutional layers on 300W Challenging subset.

Method	MobileFAN
w/o distillation	5.62
+ FA ₃	5.51
+ FA ₃ + FA ₂	5.45
+ FA ₃ + FA ₂ + FA ₁	5.40

Table 8: Mean error (%) of feature-similarity distillation of different deconvolutional layers on 300W Challenging subset.

Method	MobileFAN
w/o distillation	5.62
+ FS ₃	5.52
+ FS ₃ + FS ₂	5.47
+ FS ₃ + FS ₂ + FS ₁	5.41

Feature-aligned distillation over layer. To investigate the effectiveness of the feature-aligned distillation in facial landmark detection, we implement the feature-aligned distillation by adopting a 1×1 convolution layer to align the feature of each pixel between teacher network \mathcal{T} and student network \mathcal{S} , such that the channel of the features is matched.

We can see from TABLE 7 that feature-aligned distillation improves the performance of our proposed “MobileFAN”. By utilizing the distillation loss generated from “FA₃”, “MobileFAN” achieve 5.51% mean error on 300W Challenging subset. Moreover, we notice that the performance can be further improved by adding more layers of distillation. It is observed from TABLE 7 that “MobileFAN + FA₃ + FA₂ + FA₁” achieves 2.00% and 0.92% relative improvement over “MobileFAN + FA₃” and “MobileFAN + FA₂ + FA₃”, respectively. Similarly, the mean error of “MobileFAN + FA₂ + FA₃” has reduced from 5.51% to 5.45% compared with “MobileFAN + FA₃”. Fig. 9 provides a straightforward comparison of MobileFAN without distillation and MobileFAN with feature-aligned distillation.

Feature-similarity distillation over layer. To explore the effectiveness of feature-similarity distillation, we evaluate the mean error with various combinations of each layer on 300W Challenging subset. As can be observed from TABLE 8, applying more layers can lead to better performance.

In particular, “MobileFAN + FS₃” outperforms “MobileFAN” without distillation by a large margin, with a relative improvement of 1.78% in mean error reduction. When one more layer of distillation is added, although the improvement is marginal, the mean error is reduced from 5.52% to 5.47%. It is not surprising that “MobileFAN (0.5) + FS₃ + FS₂ + FS₁” achieves the best performance of 5.41% mean error among all versions of “MobileFAN” on 300W Challenging subset. We can see a straightforward comparison in

¹Model available from <https://github.com/wywu/LAB>.

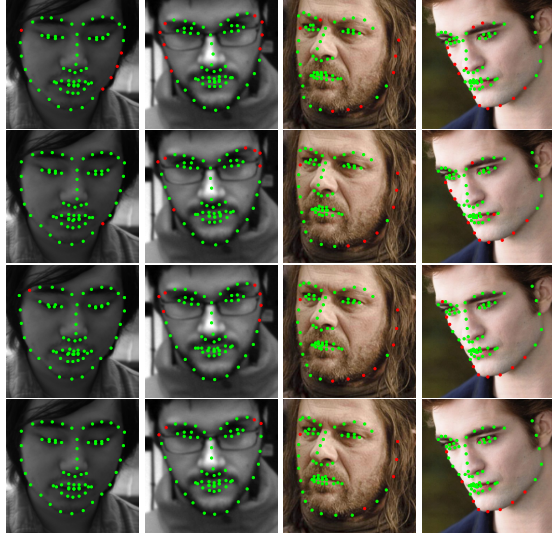


Figure 9: Comparison of MobileFAN with different distillation method. Top row: Results of MobileFAN without distillation. Second row: Results of MobileFAN with feature-aligned distillation. Third row: Results of MobileFAN with feature-similarity distillation. Bottom row: Results of MobileFAN with both feature-aligned distillation and feature-similarity distillation. Red points indicate the normalized error is larger than 0.1.

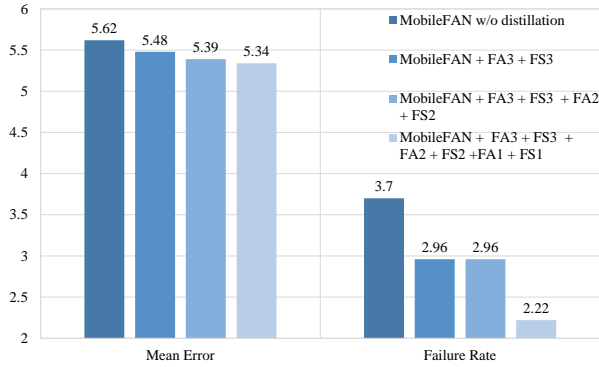


Figure 10: Normalized mean error (%) and failure rate (%) on 300W Challenging subset of various different combination of feature-aligned distillation and feature-similarity distillation.

Fig. 9 that MobileFAN with feature-similarity distillation performs better than MobileFAN without distillation.

Combination of feature-aligned distillation and feature-similarity distillation over layer. To evaluate the effectiveness of both using feature-aligned distillation and feature-similarity distillation, we report the normalized mean error and failure rate of various different combination of feature-aligned distillation and feature-similarity distillation in Fig. 10. It is shown that our “MobileFAN” with knowledge distillation performs better when more layers of both feature-aligned distillation and feature-similarity distillation are used. In particular, “MobileFAN + FA₃ + FS₃ + FA₂ + FS₂ + FA₁ + FS₁” performs better than “MobileFAN + FA₃ + FS₃” and “MobileFAN + FA₃ + FS₃ + FA₂ + FS₂”, with a relative improvement of 2.55% and 0.93% in mean error reduction, respectively. we can also find that the failure rate

of “MobileFAN + FA₃ + FS₃” is lower than that of “MobileFAN” without distillation. Similarly, the failure rate of “MobileFAN + FA₃ + FS₃ + FA₂ + FS₂ + FA₁ + FS₁” drops from 2.96% to 2.22% compared with that of “MobileFAN + FA₃ + FS₃ + FA₂ + FS₂”.

To summarize, the combination of feature-aligned distillation and feature-similarity distillation improves the performance of the compact network in facial landmark detection.

The impact of λ . To investigate the impact of λ on the training process of the proposed MobileFAN, we performed an ablation study of λ regarding mean error on COFW dataset. Experimental results in Table 9 listed the mean errors obtained with different λ (λ increases from 10^{-5} to 10^{-1}). $\lambda = 0$ means that the experiment is conducted on the proposed MobileFAN without distillation. It can be observed that the proposed MobileFAN achieved consistently lower mean errors with λ varying from 10^{-5} to 10^{-2} compared with the mean error obtained with $\lambda = 0$. The best performance is achieved at $\lambda = 10^{-2}$ on COFW dataset. Nevertheless, if λ is too large (e.g. $\lambda = 10^{-1}$), the contribution from the supervision of the ground-truth heatmaps would be limited, and thus the proposed MobileFAN may fail to converge. The overall influence of λ is positive to the training process of the proposed MobileFAN, indicating the effectiveness of the proposed distillation schemes.

Table 9: The influence of λ on COFW dataset.

λ	0	10^{-5}	10^{-4}	10^{-3}	10^{-2}	10^{-1}
Mean Error	3.82	3.78	3.79	3.74	3.66	4.06

5. Conclusion

In this paper, we focus on building a small facial landmark detection model, which remains an unsolved research problem. We propose a simple and lightweight *Mobile Face Alignment Network* (MobileFAN) by using MobileNetV2 as the encoder and three simple deconvolutional layers as the decoder. This simple design significantly helps to reduce the computational burden. With 11.5 times fewer parameters compared with the state-of-the-art models, our MobileFAN still achieves comparable or even better performance on three challenging facial landmark detection datasets. A knowledge transfer technique is proposed to enhance the performance of MobileFAN. By transferring the finer structural information encoded by the teacher Network, the performance of the proposed MobileFAN is further improved in effectiveness for facial landmark detection.

References

- [1] Xavier P Burgos-Artizzu, Pietro Perona, and Piotr Dollár. Robust face landmark estimation under occlusion. In *ICCV*, pages 1513–1520, 2013.

- [2] Xudong Cao, Yichen Wei, Fang Wen, and Jian Sun. Face alignment by explicit shape regression. In *CVPR*, pages 2887–2894, 2012.
- [3] Yu Chen, Chunhua Shen, Hao Chen, Xiu-Shen Wei, Lingqiao Liu, and Jian Yang. Adversarial learning of structure-aware fully convolutional networks for landmark localization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [4] Xuanyi Dong, Yan Yan, Wanli Ouyang, and Yi Yang. Style aggregated network for facial landmark detection. In *CVPR*, pages 379–388, 2018.
- [5] Xuanyi Dong, Shou-I Yu, Xinshuo Weng, Shih-En Wei, Yi Yang, and Yaser Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *CVPR*, pages 360–368, 2018.
- [6] Zhen-Hua Feng, Guosheng Hu, Josef Kittler, William J. Christmas, and Xiaojun Wu. Cascaded collaborative regression for robust facial landmark detection trained using a mixture of synthetic and real images with dynamic weighting. *IEEE Trans. Image Process.*, 24(11):3425–3440, 2015.
- [7] Zhen-Hua Feng, Josef Kittler, Muhammad Awais, Patrik Huber, and Xiao-Jun Wu. Wing loss for robust facial landmark localisation with convolutional neural networks. In *CVPR*, pages 2235–2245, 2018.
- [8] Zhen-Hua Feng, Josef Kittler, William J. Christmas, Patrik Huber, and Xiaojun Wu. Dynamic attention-controlled cascaded shape regression exploiting training data augmentation and fuzzy-set sample weighting. In *CVPR*, pages 3681–3690, 2017.
- [9] Golnaz Ghiasi and Charless C Fowlkes. Occlusion coherence: Localizing occluded faces with a hierarchical deformable part model. In *CVPR*, pages 2385–2392, 2014.
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [11] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [12] Sina Honari, Jason Yosinski, Pascal Vincent, and Christopher J. Pal. Recombinator networks: Learning coarse-to-fine feature aggregation. In *CVPR*, pages 5743–5752, 2016.
- [13] Chaoqun Hong, Jun Yu, Jian Wan, Dacheng Tao, and Meng Wang. Multimodal deep autoencoder for human pose recovery. *IEEE Trans. Image Process.*, 24(12):5659–5670, 2015.
- [14] Chaoqun Hong, Jun Yu, Jian Zhang, Xiongnan Jin, and Kyong-Ho Lee. Multi-modal face pose estimation with multi-task manifold deep learning. *IEEE Trans. Ind. Inf.*, 15(7):3952–3961, 2019.
- [15] Xin Jin and Xiaoyang Tan. Face alignment by robust discriminative hough voting. *Pattern Recognit.*, 60:318–333, 2016.
- [16] Marek Kowalski, Jacek Naruniec, and Tomasz Trzcinski. Deep alignment network: A convolutional neural network for robust face alignment. In *CVPRW*, pages 88–97, 2017.
- [17] Amit Kumar and Rama Chellappa. Disentangling 3d pose in a dendritic CNN for unconstrained 2d face alignment. In *CVPR*, pages 430–439, 2018.
- [18] Quanquan Li, Shengying Jin, and Junjie Yan. Mimicking very efficient network for object detection. In *CVPR*, pages 6356–6364, 2017.
- [19] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured knowledge distillation for semantic segmentation. *arXiv preprint arXiv:1903.04197*, 2019.
- [20] Xin Miao, Xiantong Zhen, Xianglong Liu, Cheng Deng, Vasilis Athitsos, and Heng Huang. Direct shape regression networks for end-to-end face alignment. In *CVPR*, pages 5040–5049, 2018.
- [21] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, pages 483–499, 2016.
- [22] Shaoqing Ren, Xudong Cao, Yichen Wei, and Jian Sun. Face alignment at 3000 fps via regressing local binary features. In *CVPR*, pages 1685–1692, 2014.
- [23] Shaoqing Ren, Xudong Cao, Yichen Wei, and Jian Sun. Face alignment via regressing local binary features. *IEEE Trans. Image Process.*, 25(3):1233–1245, 2016.
- [24] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [25] Christos Sagonas, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *ICCVW*, pages 397–403, 2013.
- [26] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, pages 4510–4520, 2018.
- [27] Sima Soltanpour, Boubakeur Boufama, and QM Jonathan Wu. A survey of local feature methods for 3d face recognition. *Pattern Recognit.*, 72:391–406, 2017.
- [28] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep convolutional network cascade for facial point detection. In *CVPR*, pages 3476–3483, 2013.
- [29] Roberto Valle, Jose M Buenaposada, Antonio Valdes, and Luis Baumela. A deeply-initialized coarse-to-fine ensemble of regression trees for face alignment. In *ECCV*, pages 585–601, 2018.
- [30] Wayne Wu, Chen Qian, Shuo Yang, Quan Wang, Yici Cai, and Qiang Zhou. Look at boundary: A boundary-aware face alignment algorithm. In *CVPR*, pages 2129–2138, 2018.
- [31] Wenyan Wu and Shuo Yang. Leveraging intra and inter-dataset variations for robust face alignment. In *CVPRW*, pages 2096–2105, 2017.
- [32] Yue Wu, Chao Gou, and Qiang Ji. Simultaneous facial landmark detection, pose and deformation estimation under facial occlusion. In *CVPR*, pages 5719–5728, 2017.
- [33] Yue Wu, Tal Hassner, KangGeon Kim, Gerard Medioni, and Prem Natarajan. Facial landmark detection with tweaked convolutional neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(12):3067–3074, 2018.
- [34] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, pages 466–481, 2018.

- [35] Shengtao Xiao, Jiashi Feng, Junliang Xing, Hanjiang Lai, Shuicheng Yan, and Ashraf Kassim. Robust facial landmark detection via recurrent attentive-refinement networks. In *ECCV*, pages 57–72, 2016.
- [36] Xuehan Xiong and Fernando De la Torre. Supervised descent method and its applications to face alignment. In *CVPR*, pages 532–539, 2013.
- [37] Jing Yang, Qingshan Liu, and Kaihua Zhang. Stacked hour-glass network for robust facial landmark localisation. In *CVPRW*, pages 2025–2033, 2017.
- [38] Jun Yu, Baopeng Zhang, Zhengzhong Kuang, Dan Lin, and Jianping Fan. iprivacy: image privacy protection by identifying sensitive objects via deep multi-task learning. *IEEE Trans. Inf. Forensics Secur.*, 12(5):1005–1016, 2016.
- [39] Jun Yu, Chaoyang Zhu, Jian Zhang, Qingming Huang, and Dacheng Tao. Spatial pyramid-enhanced netvlad with weighted triplet loss for place recognition. *IEEE Trans. Neural Netw. Learn. Syst.*, pages 1–14, 2019.
- [40] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016.
- [41] Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Occlusion-free face alignment: Deep regression networks coupled with de-corrupt autoencoders. In *CVPR*, pages 3428–3437, 2016.
- [42] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Learning deep representation for face alignment with auxiliary attributes. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(5):918–930, 2016.
- [43] Zhanpeng Zhang, Wei Zhang, Huijun Ding, Jianzhuang Liu, and Xiaoou Tang. Hierarchical facial landmark localization via cascaded random binary patterns. *Pattern Recognit.*, 48:1277–1288, 2015.
- [44] Shizhan Zhu, Cheng Li, Chen Change Loy, and Xiaoou Tang. Face alignment by coarse-to-fine shape searching. In *CVPR*, pages 4998–5006, 2015.
- [45] Ni Zhuang, Yan Yan, Si Chen, Hanzi Wang, and Chunhua Shen. Multi-label learning based deep transfer neural network for facial attribute classification. *Pattern Recognit.*, 80:225–240, 2018.