

Face Shape-Guided Deep Feature Alignment for Face Recognition Robust to Face Misalignment

Hyung-Il Kim, *Member, IEEE*, Kimin Yun, and Yong Man Ro, *Senior Member, IEEE*,

Abstract—For the past decades, face recognition (FR) has been actively studied in computer vision and pattern recognition society. Recently, due to the advances in deep learning, the FR technology shows high performance for most of the benchmark datasets. However, when the FR algorithm is applied to a real-world scenario, the performance has been known to be still unsatisfactory. This is mainly attributed to the mismatch between training and testing sets. Among such mismatches, face misalignment between training and testing faces is one of the factors that hinder successful FR. To address this limitation, we propose a face shape-guided deep feature alignment framework for FR robust to the face misalignment. Based on a face shape prior (e.g., face keypoints), we train the proposed deep network by introducing alignment processes, i.e., pixel and feature alignments, between well-aligned and misaligned face images. Through the pixel alignment process that decodes the aggregated feature extracted from a face image and face shape prior, we add the auxiliary task to reconstruct the well-aligned face image. Since the aggregated features are linked to the face feature extraction network as a guide via the feature alignment process, we train the robust face feature to the face misalignment. Even if the face shape estimation is required in the training stage, the additional face alignment process, which is usually incorporated in the conventional FR pipeline, is not necessarily needed in the testing phase. Through the comparative experiments, we validate the effectiveness of the proposed method for the face misalignment with the FR datasets.

Index Terms—Face recognition, face alignment, multi-task learning, face alignment learning, face shape prior.



1 INTRODUCTION

THANKS to the success of deep image classification and the availability of large-scale face image datasets, a deep face recognition (FR) has been actively studied. Recently, the FR performance has been considerably improved for many benchmark datasets. In particular, the performance for the LFW dataset [1] has shown about 99% accuracy in face verification [2], [3], [4]. However, when these FR algorithms are applied to the wild environment, the performance has been known to be highly degraded [5], [6], [7]. This is because benchmark datasets are usually collected from celebrities' face images with high quality, while face images in the wild environment suffer from the degradation of image quality (e.g., low-resolution, pose, and illumination variations). It leads to the mismatch of data distributions between the training and testing sets [8], [9], [10], [11], [12].

Among the mismatches, a face misalignment problem is one of the factors that hinder successful FR, which has been discussed in recent studies [9], [10], [13], [14], [15]. For example, the face misalignment problem occurs with the testing face images not elaborately aligned or differently aligned with the training face images. In order to deal with the face misalignment problem, most deep FR algorithms require the face alignment process based on the pre-defined canonical face location. For example, the recent deep FR algorithms [2], [3], [4], [16] align face images by using fiducial

points obtained from the multi-task convolutional neural network (MTCNN) [17] face detector. Furthermore, deep face keypoint estimation algorithms [18], [19], [20] have been proposed for the elaborate face alignment. However, the face keypoint estimation algorithm requires additional computational costs [21]. In addition, the face detector or additional face keypoint estimation algorithms have inherent estimation errors, which cause the face misalignment problem. More recently, to alleviate the issues related to the face alignment, there have been researches to learn a face alignment as well as a FR in an end-to-end manner (so-called *face alignment learning*) [9], [10], [13], [14], [15]. By learning the face alignment with the FR simultaneously, these works have been validated to be effective under the face misalignment since a face image is automatically aligned to the proper alignment type in testing. Motivated by the face alignment learning, in this paper, we improve the performance to be more robust to the face misalignment by learning the features being aware of face shape as well as a face image.

In this paper, we propose a face shape-guided deep feature alignment framework to address the face misalignment problem. Through both face images and the corresponding face shape priors (i.e., face keypoints), the face shape-guided feature is learned through pixel and feature alignment processes. In detail, the aggregated feature from face images and face keypoints is utilized as guidance to align two features: feature considering only face image and feature considering face shape prior as well as face image (i.e., feature alignment). Then, the proposed deep network is collaboratively trained based on three tasks (i.e., face classification, pixel and feature alignments). This shape-guided feature enables robust FR in test time without face

- H.-I. Kim and K. Yun are with Visual Intelligence Research Section, Artificial Intelligence Research Laboratory, Electronics and Telecommunications Research Institute (ETRI), Daejeon, 34129, South Korea (e-mail: {hikim, kimin.yun}@etri.re.kr)
- Y. M. Ro is with Image and Video Systems Lab, School of Electrical Engineering, Korea Advanced Institute of Science and Technology (KAIST), Daejeon, 34141, South Korea (e-mail: ymro@kaist.ac.kr)

Corresponding Author: Yong Man Ro (email: ymro@kaist.ac.kr)

alignment including keypoint estimation. The contributions of the paper can be summarized as follows:

- In training, by decoding the aggregated feature based on a face image and face shape prior, the proposed method learns the features for the well-aligned face image (*i.e.*, pixel alignment). Through the feature alignment, the face feature extraction network as an input of only face image can learn the face shape-guided feature.
- In testing, a face feature vector invariant to face alignment is extracted only from a face image based on the trained network. Since our method does not require the explicit face alignment process, we can efficiently compute the robust feature to the face misalignment.

Through the experiments, we verify the effectiveness of the proposed method conditioned under the face misalignment with face benchmark datasets.

In Section 2, we briefly discuss the previous works. Section 3 describes the proposed method. Then, the details of the experiments are presented in Section 4. Finally, we conclude the paper in Section 5.

2 RELATED WORK

2.1 Deep Face Recognition

Recently, there have been many deep FR algorithms to effectively learn discriminative features using a convolutional neural network (CNN) since the DeepFace [22]. The FaceNet [23] was proposed to learn a Euclidean space embedding by introducing a triplet loss with 200 million face images. The authors in [16] enhanced the discriminative power of the deeply learned features based on the center loss that penalizes the distances between the deep features and their corresponding class centers [16]. In order to learn angularly discriminative features, the SphereFace [4] was proposed by introducing the angular SoftMax function (*i.e.*, A-SoftMax) which imposes discriminative constraints on a hypersphere manifold [4]. Also, the CosFace [3] with a large margin cosine loss was proposed to maximize further the decision margin in the angular space [3]. In [2], the additive angular margin loss so-called ArcFace [2] was designed to obtain highly discriminative features for FR as well. More recently, to further improve FR performance, loss functions [24], [25], [26] and sample distribution-aware learning strategies [27], [28] to enhance feature discrimination are being discussed. Basically, all of these deep FR algorithms require the face alignment process to align a face image to a canonical view based on the facial keypoints.

2.2 Facial Keypoint Estimation

As discussed earlier, the face alignment process based on the facial keypoint estimation is essential to learning the deep FR network, which has been actively studied. The MTCNN [17] has been proposed to jointly learn face keypoints as well as a face bounding box, where five facial keypoints (left/right eyes, nose, left/right mouse tips) are detected with face detection. Thanks to efficient computations, the MTCNN has been widely adopted for face image-based applications. To understand a face image more accurately, the face alignment network (FAN) [18] was proposed for estimating 68 facial keypoints based on the

stacked hourglass network [29] that is widely used for human joint estimation. Furthermore, the FAN algorithm was extended to the 3D FAN algorithm for estimating 3D facial keypoints [18]. Recently, the Super-FAN [19] was proposed for estimating keypoints in the low-resolution face image, and the improved facial keypoint detector [20] robust to occlusion was proposed. Despite the improved keypoint estimation performance, the keypoint estimation network is independently trained with the FR network, and it requires additional computational complexity [21]. In addition, the facial keypoint estimation error caused by a low-quality face image promotes the vulnerability of the deep FR network to a face misalignment.

2.3 Face Alignment Learning

To deal with the aforementioned limitations, recent studies were conducted to simultaneously train a face localization network for face alignment and a face feature extraction network for FR, which is called Face Alignment Learning. Zhong *et al.* proposed the end-to-end learning framework [13] for estimating the face image’s transformation and FR by the Spatial Transformer Network (STN) [30]. Here, STN’s localization network predicted the 2D transform parameters of the face image to transform the face image into a canonical view. However, the estimated transformation can capture only coarse geometric information as a holistic parametric model [10]. To deal with the accurate transformation (*i.e.*, non-rigid transformation), the recursive spatial transformer (ReST) [14] by progressively aligning face images was proposed. Despite the effective approach, the progressive aligning pipeline causes the degradation of the face image according to the repetitive rectifications for an image and features. Besides, GridFace [10] was proposed to reduce geometric facial variations and improve the recognition performance by rectifying the face by local grid-level homography transformations [10]. Recently, an adaptive pose alignment (APA) [9] has been proposed for aligning each face of training or test set to optimal alignment templates according to the facial pose instead of a predefined template [9]. In [15], Wei *et al.* proposed the adaptive alignment of the face feature map based on the warp grid obtained from the localization network trained with face keypoints. In this paper, we propose the end-to-end face feature learning framework guided by the face shape as the face alignment learning instead of a direct transformation of a face image or feature maps.

3 PROPOSED METHOD

In this paper, we propose a face shape-guided deep feature alignment framework robust to face misalignment. By using both a randomly cropped face image (\mathbf{x}_i^r) from the i -th training face image (\mathbf{x}_i) and the corresponding face shape prior, the face shape-guided feature is trained based on the well-aligned face images (\mathbf{x}_i^w). Note that the well-aligned face images denote face images cropped by considering face keypoints. And, the randomly cropped face images are obtained by over-sampling strategy [31]. Please refer to Preprocessing in Section 4.1 for details. When testing, our face shape-guided feature vector, which is invariant

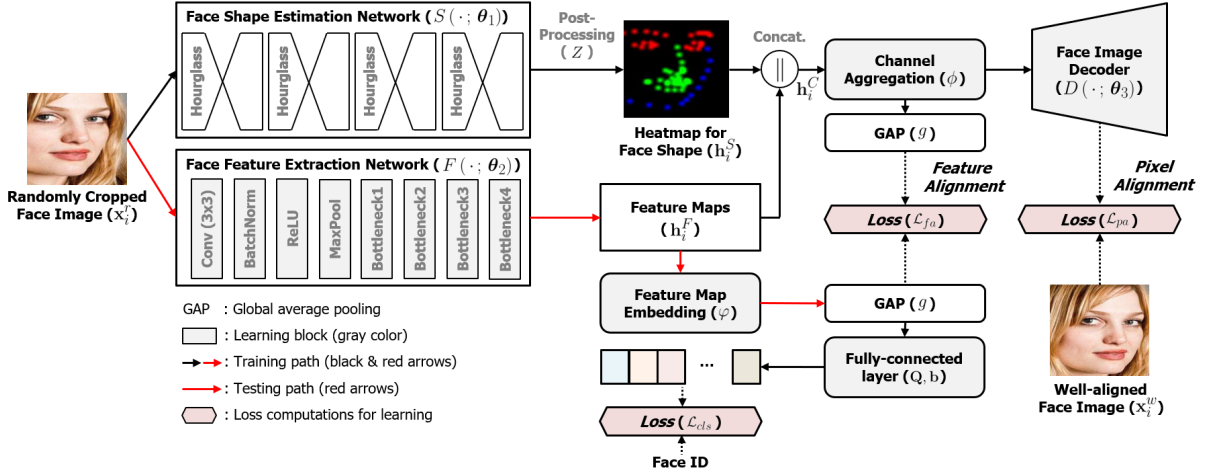


Fig. 1. Overview of the proposed face shape-guided deep feature alignment framework. The proposed method is mainly comprised of the face shape estimation network, face feature extraction network, and face image decoder. Based on three networks, two alignment processes (*i.e.*, pixel alignment and feature alignment) are introduced. Then, through the pixel alignment and feature alignment processes, the face shape-guided feature is trained by jointly classifying the face and reconstructing the well-aligned face image. Note that only the face feature extraction network is used for the testing phase. Best viewed in color.

to the face alignment, enables us to achieve robust FR to face misalignment without the additional face alignment process. In the following subsections, we describe more details of the proposed method.

3.1 Architecture

As shown in Fig. 1, the proposed method mainly consists of three parts for the purpose of face feature extraction robust to face misalignment: 1) Face shape estimation network (S), 2) Face feature extraction network (F), and 3) Face image decoder (D), which are connected by two separated neural layers (*i.e.*, feature map embedding (φ) and channel aggregation (ϕ)). The face shape estimation network infers the coordinates of the facial keypoints as a face shape prior, resulting in the face shape’s heatmaps. And, the face feature extraction network outputs face appearance features based on a backbone network. After aggregating outputs from the face shape estimation and feature extraction networks, the well-aligned face image is reconstructed by the face image decoder from the randomly cropped face image (*i.e.*, pixel alignment). Through the pixel alignment, we can model the characteristic of the well-aligned face image. The face shape-aggregated feature is simultaneously connected to the face feature extraction network through the feature alignment. It helps our face feature extraction network to learn the face shape-guided feature. In other words, the feature alignment is helpful for extracting the face shape-guided feature without the help of the face shape estimation network in testing. By learning the characteristic of the well-aligned face image in an end-to-end manner, we can extract face features robust to face misalignment in the testing phase without an explicit face alignment process.

3.1.1 Face Shape Estimation Network

In order to estimate the face shape prior, we use the 2D face alignment network (FAN) [18], which consists of four consecutive hourglass modules [29]. In this paper, the FAN architecture and the corresponding pre-trained model are

utilized for the face shape estimation without any modification, where the pre-trained parameters are frozen in training to obtain the face shape prior stably. For the input face image (x_i^r , *i.e.*, randomly cropped face image), the network parameterized by θ_1 infers 68-channel heatmaps for the corresponding 68 facial keypoints. To take advantage of the heatmaps effectively, we perform the following post-processing method: Gaussian blurring, resizing, and channel conversion. First, the 68-channel heatmaps are blurred by a Gaussian kernel with σ for the location of the heatmap’s peak point (*i.e.*, the facial keypoint location). This is to emphasize the importance of the surrounding areas around the keypoints as suggested in [32]. Then, we resize the heatmaps to 56×56 to align with the feature maps extracted from the face feature extraction network. Note that the 68-channel heatmaps for the estimated keypoints are converted to 3-channel heatmaps for efficient memory consumption. In other words, the converted heatmap is made to have three image channels, where the channel corresponds to each part: 1) R-channel: heatmap for eyes and eyebrows, 2) G-channel: heatmap for the nose and mouth, and 3) B-channel: heatmap for the face boundary as shown in Fig. 1. In summary, the extracted heatmap (h_i^S) is represented as follows:

$$h_i^S = Z[S(x_i^r; \theta_1)], \quad (1)$$

where S denotes the operation of the FAN’s feedforward computation. And, Z means the post-processing function for the Gaussian blurring, resizing, and channel conversion.

3.1.2 Face Feature Extraction Network

To extract face appearance features, we design the face feature extraction network (F) by modifying the ResNet50 [33]. Like the original ResNet50, the F is comprised of 2D convolution, batch normalization, ReLU activation, max-pooling layers, and four Bottleneck layers [33] in order. Here, the number of channels in each layer and the stride in the Bottleneck layers are modified to align feature maps with the heatmaps for the face shape prior. Specifically, the number of channels is all halved, *i.e.*, 512 channels are changed to

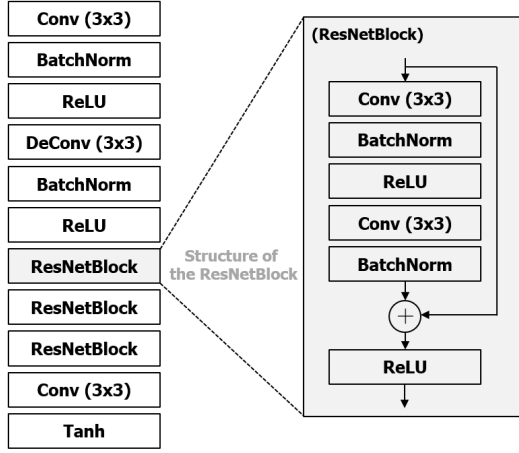


Fig. 2. The structure of the face image decoder (D), where the structure of the ResNetBlock (i.e., Basic block [33] in the D is presented in the gray box.). Note that the 'Conv' and 'DeConv' denote convolution and deconvolution operations, respectively.

256 channels in the fourth Bottleneck layer. And, the stride for the bottleneck layer is set to 1 instead of 2. For the \mathbf{x}_i^r , the face appearance feature maps (\mathbf{h}_i^F) are obtained by the following equation:

$$\mathbf{h}_i^F = F(\mathbf{x}_i^r; \theta_2), \quad (2)$$

where θ_2 is the parameters of the face feature extraction network.

3.1.3 Face Image Decoder

The face image decoder is incorporated to learn the characteristic of the well-aligned face image \mathbf{x}_i^w by using features from \mathbf{h}_i^S and \mathbf{h}_i^F . Prior to the decoding, both the \mathbf{h}_i^S and \mathbf{h}_i^F are concatenated to form a stacked feature map (\mathbf{h}_i^C) as:

$$\mathbf{h}_i^C = [\mathbf{h}_i^S \parallel \mathbf{h}_i^F], \quad (3)$$

where " \parallel " symbol means the concatenation operation in a channel direction. Then, the stacked feature \mathbf{h}_i^C is effectively aggregated to make a face shape-guided feature through the Channel Aggregation. The channel aggregation layer is simply comprised of one 1×1 convolution layer and one batch normalization layer, which fuses face appearance features and face shape prior in a channel-wise manner. For simplicity, we denote the channel aggregation operation as ϕ . The input of the face image decoder can be represented as $\phi(\mathbf{h}_i^C)$. The decoded result is obtained by the following equation:

$$\tilde{\mathbf{x}}_i^w = D(\phi(\mathbf{h}_i^C); \theta_3), \quad (4)$$

where the face image decoder (D) is parameterized by θ_3 . The decoder network firstly reduces the number of concatenated feature maps to 64 by a 3×3 convolution layer. The next 3×3 deconvolution layer is used for up-sampling the feature map to double the resolution. Then, three residual blocks are used to decode features. Finally, a 3×3 convolution layer reconstructs the well-aligned face image. The detailed architecture of the face image decoder is shown in Fig. 2. Note that the process to reconstruct the well-aligned face image from the stacked features of the randomly cropped face image is called to *Pixel Alignment*.

In addition to the procedure mentioned above, an additional learning path should be devised to extract face features that are robust to face misalignment in the testing. Motivated by the concept of the feature alignment in the recent researches [34], [35], we design the learning path for the *Feature Alignment*. In other words, the feature alignment in our network enables us to train the face feature extraction network guided by the channel aggregated feature used for the face image decoder. Since the channel aggregated feature includes information related to the face appearance and the face shape, the additional learner φ learns the function to map the face appearance feature \mathbf{h}_i^F extracted from the face feature extraction network F to the guidance ($\phi(\mathbf{h}_i^C)$). Note that we transform both feature maps $\phi(\mathbf{h}_i^C)$ and $\varphi(\mathbf{h}_i^F)$ to feature vectors by the global average pooling (GAP) for efficient computations [36]. In Table 1, the details of the architecture to be trained in our framework are summarized.

Algorithm 1: Pseudo code for training the proposed method. All training samples are divided into n_B batches and used for training.

Input: Randomly cropped and well-aligned face image pairs and corresponding ID labels, Parameters θ_1 for Face Shape Estimation Network, Number of epochs n_E , Number of batches n_B , Learning rate δ , Momentum τ

Output: $\Theta = \{\theta_2, \theta_3, \phi, \varphi, \mathbf{Q}, \mathbf{b}\}$

```

1 for  $t = 1, \dots, n_E$  do
2   for  $b = 1, \dots, n_B$  do
3     # Feed-forward Operation
4     Forward propagating  $S$  by  $\theta_1$  and post
      -processing  $Z$  to obtain face shape priors;
5     Forward propagating  $F$  by  $\theta_2$  to obtain the
      face feature maps;
6     Concatenating the face shape prior and face
      feature maps, then forward propagating  $\phi$ ;
7     Forward propagating  $\varphi, \mathbf{Q}, \mathbf{b}$  to obtain the
      estimated ID labels;
8     Forward propagating  $D$  by  $\theta_3$  to obtain the
      reconstructed face images;
9
10    # Loss Computation
11    Compute  $\mathcal{L}_{cls}$  by Eq. (5);
12    Compute  $\mathcal{L}_{pa}$  by Eq. (6);
13    Compute  $\mathcal{L}_{fa}$  by Eq. (7);
14    Weighted sum ( $\mathcal{L}$ ) of loss functions by Eq. (8);
15
16    # Parameter Update
17    if  $b = 1$  then
18       $\mathbf{v}_b^t \leftarrow \nabla_{\Theta} \mathcal{L}$ ;
19    else
20       $\mathbf{v}_b^t \leftarrow \tau \mathbf{v}_{b-1}^t + \nabla_{\Theta} \mathcal{L}$ ;
21    Update:  $\Theta_{b+1}^t \leftarrow \Theta_b^t - \delta^t \mathbf{v}_b^t$ ;
22   $\Theta_1^{t+1} \leftarrow \Theta_{n_B}^t$ 

```

3.2 Training

In order to train the proposed network, three-loss functions are introduced. First, a cross-entropy loss is used to classify a

TABLE 1

Summary of the proposed deep network's architecture. $[\cdot]$ represents the operations of the ResNet blocks (Bottleneck [33] and Basic [33]). DeConv layer means the deconvolution operation to enlarge the resolution of a feature map. And, 8631 in the fully connected layer stands for the number of classes in the training dataset.

Network	Layer Name	Output Size	Operation
Face Feature Extraction Network (F)	Conv1	112×112	$7 \times 7, 64$, stride 2 (BatchNorm, ReLU)
	Conv2_x (Bottleneck)	56×56	3×3 Max Pooling, stride 2
			$1 \times 1, 32$, stride 1
			$3 \times 3, 32$, stride 3
	Conv3_x (Bottleneck)	56×56	$1 \times 1, 128$, stride 1
Face Image Decoder (D)	Conv2_x (Bottleneck)	56×56	$1 \times 1, 64$, stride 1
	Conv3_x (Bottleneck)	56×56	$3 \times 3, 64$, stride 3
			$1 \times 1, 256$, stride 1
			$1 \times 1, 128$, stride 1
	Conv4_x (Bottleneck)	56×56	$3 \times 3, 128$, stride 3
Feature Map Embedding (φ)	Conv1	56×56	$1 \times 1, 512$, stride 1 (BatchNorm, ReLU)
	DeConv	112×112	$3 \times 3, 64$, stride 2 (BatchNorm, ReLU)
	Conv2_x (Basic)	112×112	$3 \times 3, 64$, stride 1
	Conv3_x (Basic)	112×112	$3 \times 3, 64$, stride 1
	Conv4_x (Basic)	112×112	$3 \times 3, 64$, stride 1
Channel Aggregation (ϕ)	Conv5	112×112	$3 \times 3, 3$, stride 1 (Tanh)
	Conv	56×56	$1 \times 1, 512$, stride 1 (BatchNorm)
			$1 \times 1, 512$, stride 1 (BatchNorm)
			$1 \times 1, 512$, stride 1 (BatchNorm)
	Fully Connected Layer (\mathbf{Q}, \mathbf{b})	1×1	Global Average Pooling, 8631 (SoftMax)

face image into one of the classes for the embedded feature $g(\varphi(\mathbf{h}_i^F))$ after the feature extraction by F and the GAP, which is defined by

$$\mathcal{L}_{cls} = -\frac{1}{nc} \sum_{\forall i} \sum_{\forall c} y_i^c \log \tilde{y}_i^c, \quad (5)$$

where $y_i^c \in \{0, 1\}$ is the c -th element of the one-hot vector corresponding to the ground truth class label of the i -th sample, and \tilde{y}_i^c is the c -th element of the estimated label by a SoftMax function, i.e., $\text{SoftMax}(\mathbf{Q}^T g(\varphi(\mathbf{h}_i^F)) + \mathbf{b})$. The \mathbf{Q} and \mathbf{b} are the weight matrix and bias vector for the fully-connected layer. The n and c are the number of samples used for training in an epoch and the number of classes, respectively. Then, to learn the face image decoder, the L1 loss between the well-aligned face image and the decoded face image, i.e., loss for the pixel alignment, is defined by

$$\mathcal{L}_{pa} = \frac{1}{n} \sum_{\forall i} \|\mathbf{x}_i^w - \tilde{\mathbf{x}}_i^w\|_1, \text{ where } \tilde{\mathbf{x}}_i^w = D(\phi(\mathbf{h}_i^C); \theta_3). \quad (6)$$

Finally, the loss for the feature alignment to learn the face shape-guided feature is defined by the following equation:

$$\mathcal{L}_{fa} = \frac{1}{n} \sum_{\forall i} \|g(\phi(\mathbf{h}_i^C)) - g(\varphi(\mathbf{h}_i^F))\|_2^2. \quad (7)$$

From Eq. (7), the feature considering both the face image and face shape prior, $g(\phi(\mathbf{h}_i^C))$, is distilled to the feature

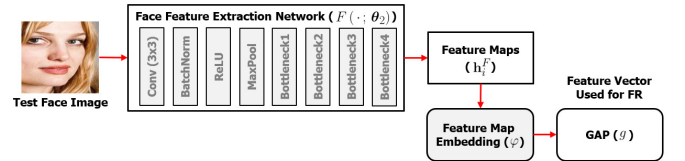


Fig. 3. Inference step for the proposed method.

from the face feature extraction network, $g(\varphi(\mathbf{h}_i^F))$. By using three loss functions, the total loss (\mathcal{L}) for training the proposed deep network is defined as:

$$\mathcal{L} = \alpha \mathcal{L}_{cls} + \beta \mathcal{L}_{pa} + \gamma \mathcal{L}_{fa}, \quad (8)$$

where α, β , and γ denote the hyper-parameters that control the balance for the total loss function. We optimize the following objective function to obtain the learning parameters (i.e., $\Theta = \{\theta_2, \theta_3, \phi, \varphi, \mathbf{Q}, \mathbf{b}\}$) by mini-batch gradient descent:

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}. \quad (9)$$

Based on three-loss functions, we train the proposed deep network with three steps: 1) feed-forward operation, 2) loss computation, and 3) parameter update by mini-batch gradient descent. The training method is summarized in Algorithm 1.

TABLE 2

Summary of data preprocessing and hyper-parameters for training the proposed network. ([†] All face images were normalized to have values ranging from -1 to 1, dividing 128.0 after subtracting 127.5 from original pixels values after the preprocessing. * The learning rate (δ) was decayed by the factor of 0.1 for every 30 epochs.)

Data preprocessing [†]			
Well-aligned face image (\mathbf{x}^w)			
- Crop original face image based on the bounding box			
- Align face image based on the keypoints via similarity transform			
- Resize the face image to 224×224 pixels			
Randomly cropped face image (\mathbf{x}^r)			
- Adjust bounding box coordinates with 10 pixels margin			
- Resize the face image to 256×256 pixels			
- Randomly crop 224×224 pixels within the 256×256 pixels face image			
Hyper-parameters for training the proposed deep network			
Learning rate* (δ)	0.1	Momentum (τ)	0.9
Number of epochs (n_E)	100	Number of batches (n_B)	64
Weight for \mathcal{L}_{cls} (α)	1.0	Weight for \mathcal{L}_{pa} (β)	1.0
Weight for \mathcal{L}_{fa} (γ)	1.0	Gaussian blurring (σ)	2.0

3.3 Inference

As can be seen in Fig. 3, given the learned face feature extraction network F , an input face image is processed to extract a feature vector (*i.e.*, $g(\varphi(\mathbf{h}^F))$) robust to the face misalignment regardless of the networks S and D . The extracted feature vector is used for the FR. Note that even if the proposed face shape-guided deep feature alignment framework requires a face shape prior extracted by S , the face shape prior is no longer required in the testing step. Thus, it can reduce the computational complexity compared to the conventional FR pipeline equipped with an explicit face alignment process.

4 EXPERIMENTS

4.1 Experimental Settings

Dataset. For the experiments, the large-scale VGGFace2 [31] dataset was used for training the proposed method. The VGGFace2 dataset is comprised of 3.31 million face images from 9,131 identities which have large variations in pose, age, illumination, ethnicity, and profession [31]. As suggested in [31], face images from 8,631 persons were used for training the proposed deep network, and the remaining face images for 500 disjoint persons were used for the validation. To evaluate the proposed method, we adopted four benchmark datasets: 1) LFW [1], 2) Cross-Age LFW (CALFW) [37], 3) Cross-Pose LFW (CPLFW) [45], 4) YTF [38], and 5) MegaFace [39]. The LFW dataset contains 13,233 face images from 5,749 face identities and provides 6,000 pairs for face verification test [1]. As the reorganized version of the LFW dataset, the CALFW dataset is comprised of 6,000 pairs with large-age variations [37] for face verification tests as well. The CPLFW dataset was constructed by searching and selecting 3,000 positive face pairs with pose differences to add pose variation to intra-class variance [45]. And 3,000 negative pairs with the same gender and race were also selected to reduce the influence of attribute difference between positive and negative pairs [45]. The YTF dataset includes 3,425 videos of 1,595 identities, where two video clips are verified whether they are matched or not [38]. For a face verification with four datasets, we

measured accuracy with 10-fold cross-validation. For a face identification, we used MegaFace dataset which consists of 1 million distractors with the Challenge 1 of 100,000 face images from 530 celebrities (*i.e.*, FaceScrub [40]). Note that the cleaned version of the MegaFace dataset was adopted in our experiment as introduced in [2]. For all experiments, the 512-dimensional feature vectors for input face images were extracted, then normalized feature vectors to have a unit norm by a L2 normalization. The distance between face images in verification pairs (LFW, CALFW, and YTF), as well as the gallery-probe list (MegaFace), was computed based on a cosine distance.

Preprocessing. As discussed earlier, two types of face images were used for training. The well-aligned face image (\mathbf{x}^w) was obtained by cropping an original face image based on the bounding box coordinates provided in the VGGFace2 dataset and keypoints by resizing the face images to 224×224 pixels. In contrast, the randomly cropped face images (\mathbf{x}^r) were obtained by adjusting bounding box coordinates with 10 pixels margin, then resized to 256×256 pixels. Finally, the face images were randomly cropped to have 224×224 pixels. Pixel intensities were normalized to have values ranging from -1 to 1 by dividing 128.0 after subtracting 127.5 from original pixel values.

Training Settings. To train the proposed deep network, four NVIDIA Titan Xp GPUs with 12GB GPU memory were used. The batch size for one epoch was set to 64 samples (*i.e.*, 64 well-aligned and 64 randomly cropped face images as a pair with the corresponding face ID labels). The learning rate (δ) and momentum (τ) parameters were set to 0.1 and 0.9 recommended in [31], [33], respectively. Specifically, since the pre-trained model for the proposed network is not available and learning from scratch is required, a relatively large learning rate was set. The learning rate for training was decayed by the factor of 0.1 for every 30 epochs. We set the total number of epochs to 100. Note that the face shape estimation network was frozen to provide stable keypoint heatmaps in training. The image size to be decoded in the face image decoder was set to 112×112 pixels to reduce the memory overhead. And, the standard deviation in the Gaussian blurring for generating heatmaps from the keypoint coordinates was set to 2 (*i.e.*, $\sigma = 2$). The hyper-parameters in Eq. (8) were all set to 1 (*i.e.*, $\alpha = \beta = \gamma = 1$). Finally, the proposed method was optimized by a mini-batch gradient descent algorithm. The data preprocessing and training settings are summarized in Table 2.

4.2 Effect of the Face Misalignment for Deep FR

In this experiment, we first evaluated the LFW dataset in order to investigate the effect of the face misalignment for deep FR algorithms by adjusting the degree of the face misalignment. To adjust the face alignment for testing face images, we introduced margin parameters $m = (m_{x_1}, m_{x_2}, m_{y_1}, m_{y_2})$, where each element denotes a ratio controlling the margin for each point of a bounding box. Here, the face image's bounding box was defined by the two points obtained from a face detector: the left top point $\mathbf{x}_1 = (x_1, y_1)^\top$ and the right bottom point $\mathbf{x}_2 = (x_2, y_2)^\top$. Given the face image's bounding box, the bounding box



Fig. 4. Example face images used for the experiments in order to investigate the effect of the face misalignment depending on the different margin parameters.

TABLE 3

Face verification performance (Accuracy, %) for the LFW dataset depending on the margin parameters, where “Optimal Alignment” means that the required optimal alignment is satisfied for each deep FR method.

	Margin parameters to control face alignment							Optimal Alignment
	m_1	m_2	m_3	m_4	m_5	m_6	m_7	
SphereFace [4]	86.82 ± 1.34	79.55 ± 2.33	75.15 ± 1.96	72.85 ± 2.35	72.37 ± 2.66	73.47 ± 2.30	75.15 ± 1.80	99.10 ± 0.38
CosFace [3]	97.25 ± 1.10	93.40 ± 1.07	85.55 ± 1.55	80.77 ± 1.76	80.20 ± 1.77	82.95 ± 1.11	82.75 ± 2.07	99.52 ± 0.30
VGGFace2 [31]	98.98 ± 0.55	95.95 ± 0.84	84.63 ± 1.86	77.85 ± 1.89	77.87 ± 1.91	83.50 ± 1.83	86.75 ± 1.26	99.08 ± 0.54
ArcFace [2]	94.72 ± 0.70	93.12 ± 1.29	84.38 ± 1.44	79.95 ± 1.99	79.03 ± 2.08	83.07 ± 1.80	76.87 ± 2.28	99.72 ± 0.19
Proposed	99.28 ± 0.58	99.30 ± 0.44	98.85 ± 0.65	97.47 ± 1.06	97.12 ± 1.10	98.80 ± 0.79	98.67 ± 0.73	99.30 ± 0.44

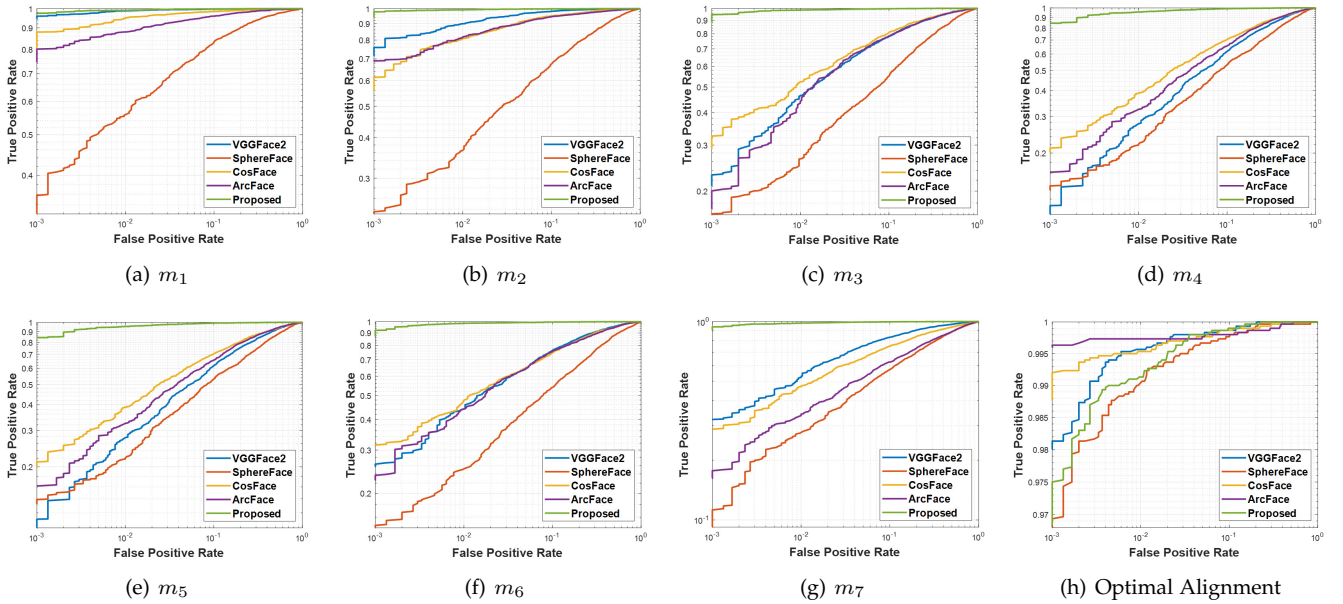


Fig. 5. Face verification ROC curves for the LFW dataset depending on the margin parameters. Best viewed in color.

for the misaligned face image (i.e., $\mathbf{x}'_1 = (x'_1, y'_1)^\top$ and $\mathbf{x}'_2 = (x'_2, y'_2)^\top$) was obtained by the following equation:

$$\begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix} = \begin{pmatrix} 1 + 0.5m_{x_1} & -0.5m_{x_1} \\ -0.5m_{x_2} & 1 + 0.5m_{x_2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad (10)$$

$$\begin{pmatrix} y'_1 \\ y'_2 \end{pmatrix} = \begin{pmatrix} 1 + 0.5m_{y_1} & -0.5m_{y_1} \\ -0.5m_{y_2} & 1 + 0.5m_{y_2} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}. \quad (11)$$

For the evaluation, we generated the seven types of misaligned face images: $m_1 = (0.50, 0.50, 0.50, 0.50)$, $m_2 = (1.00, 1.00, 1.00, 1.00)$, $m_3 = (1.50, 1.50, 1.50, 1.50)$, $m_4 = (2.00, 2.00, 2.00, 2.00)$, $m_5 = (2.50, 2.50, 2.50, 2.50)$, $m_6 = (1.25, 0.70, 1.75, 2.15)$, and $m_7 = (0.33, 2.13, 2.17, 2.34)$. The margin parameters from m_1 to m_5 adjusted the face bounding boxes with the same ratio for all directions (left/right/top/bottom). The m_6 and m_7 adjusted the face bounding boxes differently for all directions, where each element was randomly selected. Fig. 4 shows the example face images depending on the margin parameters used in

this experiment. By increasing the margin parameter, we obtain the face image similar to the original input face image before cropping. For the evaluation of the robustness to the face misalignment, the state-of-the-art deep FR algorithms were compared: 1) SphereFace [4] (64-layer CNN with a residual unit trained with CASIA WebFace [41]), 2) CosFace [3] (64-layer CNN with a residual unit trained with CASIA WebFace), 3) VGGFace2 [31] (50-layer CNN with a residual unit trained with VGGFace2 dataset [31]), and 4) ArcFace [2] (101-layer CNN with an improved residual unit trained with MS-Celeb-1M [42]). Note that the SphereFace, CosFace, and ArcFace used basically horizontal flipping as data augmentation for training the deep networks [2], [3], [4]. In contrast, the VGGFace2 algorithm adopted the over-sampling strategy [31] that extends the face bounding box, then randomly crops as the data augmentation. Based on the publicly available deep FR models, we measured the accuracy in our machine under the same conditions.

Table 3 shows the face verification performance for the

TABLE 4
Face recognition accuracy (%) for the LFW and CALFW datasets under the different types of alignment conditions.

	LFW			CALFW		
	Optimal Alignment	Random Alignment	Without Detection	Optimal Alignment	Random Alignment	Without Detection
SphereFace [4]	99.10 \pm 0.38	56.87 \pm 2.03	72.27 \pm 2.52	89.53 \pm 0.52	52.78 \pm 2.33	56.72 \pm 2.47
CosFace [3]	99.52 \pm 0.30	64.60 \pm 1.28	80.12 \pm 1.67	90.62 \pm 0.41	58.32 \pm 1.56	61.83 \pm 1.63
VGGFace2 [31]	99.08 \pm 0.54	85.71 \pm 1.85	77.87 \pm 1.82	87.75 \pm 0.70	63.80 \pm 1.75	59.57 \pm 1.81
ArcFace [2]	99.72\pm 0.19	64.09 \pm 1.81	79.00 \pm 2.03	93.63\pm 0.22	59.93 \pm 1.80	63.85 \pm 1.94
Proposed	99.30 \pm 0.44	97.46\pm 0.95	97.17\pm 0.98	90.61 \pm 0.51	89.14\pm 0.82	88.45\pm 0.88

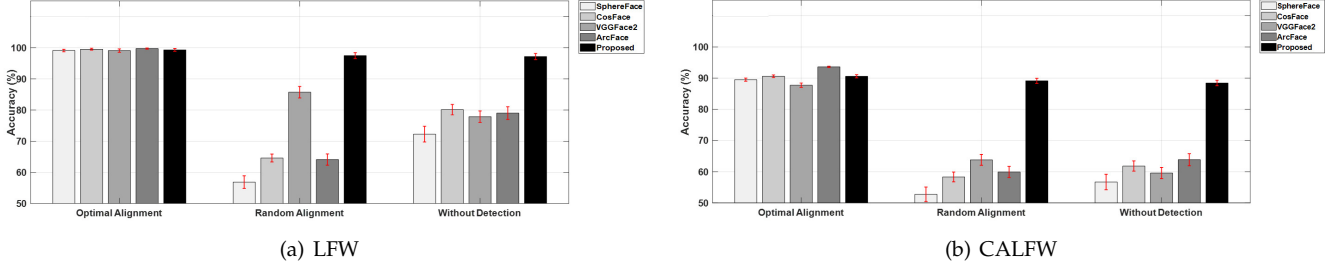


Fig. 6. Error bar graphs for the face recognition accuracy (%) for the LFW and CALFW datasets presented in Table 4.

LFW dataset depending on the margin parameters. And Fig. 5 shows the face verification receiver operating characteristic (ROC) curves for the LFW dataset depending on the margin parameters. The performance (*i.e.*, average performance and standard deviation) was measured based on the splits for 10-fold cross-validation in [1]. The standard deviation can be thought of as a measure of how stable a recognition performance is for repeated experiments, and a small standard deviation value is preferred. Note that the “Optimal Alignment” means the face verification results accompanied by the optimal alignment process used in each deep FR algorithm. For the three deep FR algorithms (SphereFace, CosFace, and ArcFace), a face image is aligned based on the facial keypoints, then crop the face images with 96×112 pixels (for SphereFace), 112×112 pixels (for both CosFace and ArcFace). The VGGFace2 algorithm crops face images based on the bounding boxes without facial keypoint estimation. Even if the previous deep FR algorithms showed the high performance for the face images obtained from the optimal face alignment, the accuracy was significantly degraded when changing the type of face alignment. It shows that the performance of the existing methods becomes highly sensitive to face alignment. Furthermore, standard deviation values become large (*i.e.*, decreased stability) as the type of face alignment changes. In the case of the SphereFace, the degradation of the performance was significant. This is mainly because the input face image of the SphereFace has a different ratio for height and width. By changing the margin parameter, the ratio of a face image is considerably degraded. In the case of the VGGFace2, it could be tolerant as the margin parameter varies thanks to the over-sampling strategy in Table 3. However, the accuracy for the VGGFace2 algorithm was also degraded according to the increase of the margin parameter despite the advantage of the over-sampling. In contrast, the proposed method showed robust performance with a relatively small standard deviation to the change of the margin parameters. Even more, we observe that the proposed method shows comparable accuracy though the face image is cropped with

high margin values including backgrounds. In other words, the basic data augmentation strategy like the over-sampling can be used as one of the solutions that resolve the face misalignment problem, however, it would not be an optimal solution.

4.3 Robustness under Random Alignment for Deep FR

In the conventional FR pipeline, there are errors and uncertainties caused by face detection and face alignment. These uncertainties of the face detector and face keypoint detector make the face images cropped with a different type of alignment. In this experiment, we changed the margin parameters randomly instead of a specific margin parameter as presented in Section 4.2 (we denote the case as “Random Alignment”). Here, the margin parameters were randomly generated from 0 to 3 by a uniform distribution (*i.e.*, $\mathcal{U}(0, 3)$), which were applied to four elements of the margin parameter independently. Additionally, we compared the accuracy under the original test image as it is (without a face detection, which is denoted as “Without Detection”) as the extreme case. Note that the accuracy of “Optimal Alignment” for the proposed method was selected by the maximum performance between m_1 and m_7 . Table 4 shows the face verification performance under three alignment conditions for the LFW and CALFW datasets with error bar graphs in Fig. 6. Without the optimal face alignment, the accuracies for the previous methods were significantly degraded as observed in Section 4.2. In particular, we observe that the accuracy for “Random Alignment” is worse than that of “Without Detection”. This is mainly caused by the mismatches between test face images in a pair. In contrast, the proposed method showed robust performance with a relatively small standard deviation even under different types of face misalignment conditions. Of course, we can observe few failure cases for the proposed method: 1) when the face bounding box contains large amounts of background information other than the face area as it is extended by the margin parameter, 2) when the partial face not including



Fig. 7. Example face images for the failure cases of the proposed method.

TABLE 5

LFW face verification accuracy (%) depending on the combination of the loss functions: \mathcal{L}_{cls} (loss for classification), \mathcal{L}_{pa} (loss for pixel alignment), and \mathcal{L}_{fa} (loss for feature alignment).

Loss functions for training	LFW	
	Random Alignment	Without Detection
\mathcal{L}_{cls}	84.01 ± 1.43	74.07 ± 1.62
$\mathcal{L}_{cls} + \mathcal{L}_{pa}$	92.65 ± 1.14	93.05 ± 1.22
$\mathcal{L}_{cls} + \mathcal{L}_{fa}$	95.32 ± 0.83	95.15 ± 0.91
$\mathcal{L}_{cls} + \mathcal{L}_{pa} + \mathcal{L}_{fa}$	97.46 ± 0.95	97.18 ± 0.98
$\mathcal{L}_{cls} + 2\mathcal{L}_{pa} + \mathcal{L}_{fa}$	97.23 ± 0.92	97.00 ± 0.87
$\mathcal{L}_{cls} + \mathcal{L}_{pa} + 2\mathcal{L}_{fa}$	97.62 ± 0.88	97.08 ± 0.94

all face components (eyes, nose, or mouth) by the random alignment is used for testing. Fig. 7 shows the example failure cases for the proposed method. These failure cases were included in order to test the proposed method under extreme conditions in the experiment. However, we can observe that the proposed method shows robustness to extreme conditions compared to the previous works.

4.4 Exploratory Experiments

In this section, we investigate the effectiveness of the proposed method via exploratory experiments: 1) effect of the loss functions, and 2) effect of the face shape estimation network.

Effect of the Loss Functions. First, we reported the accuracy depending on the loss functions to investigate the effect of the three different loss functions. According to Table 5, when the classification loss (\mathcal{L}_{cls}) was only considered, the performance was sensitive to the face misalignment and showed similar accuracy to VGGFace2. In contrast, when the pixel alignment loss (\mathcal{L}_{pa}) reflecting the reconstruction loss between the well-aligned and randomly cropped face images was additionally used, the accuracy was improved. We observe that the additional shape-related information by the pixel alignment loss gives a positive effect to train the deep network even though there is no direct guidance from the feature alignment process. Additionally, when the feature alignment loss (\mathcal{L}_{fa}) was considered with the \mathcal{L}_{cls} , the accuracy was more enhanced thanks to the direct guidance of the shape-aggregated feature. Finally, we achieved the best performance when all loss functions were accommodated. Also, we evaluated the accuracy for the three-loss functions with different hyper-parameters that control the balance for the total loss function. As can be seen in Table 5, the additional performance gain could be obtained when the feature alignment loss was more weighted than the pixel alignment loss. It turns out that the feature alignment loss is much more important than the pixel alignment loss. In

TABLE 6

LFW face verification accuracy (%) depending on the face shape estimation network, where the keypoint estimation performance (%) by the face shape estimation network was measured AUC @ 8% NME [43] with 300W private dataset [44].

Face Shape Estimation Network (FSEN)	Keypoint Estimation Performance	LFW (1:1 Verification)	
		Optimal Alignment	Random Alignment
Proposed w/o FSEN	-	99.04 ± 0.71	95.25 ± 0.60
Proposed w/ FSEN (FAN $_{\epsilon=4}$)	38.34	98.97 ± 0.54	96.78 ± 0.56
Proposed w/ FSEN (FAN $_{\epsilon=2}$)	45.49	99.23 ± 0.65	97.46 ± 0.58
Proposed w/ FSEN (FAN $_{\epsilon=0}$)	48.38	99.30 ± 0.44	97.46 ± 0.95

conclusion, this experiment demonstrates that the three-loss functions are complementary to each other even if the accuracy can vary by the hyper-parameters selection. Combining all functions is highly effective for the purpose of FR robust to face misalignment.

Effect of the Face Shape Estimation Network. In order to investigate the contribution of the face shape information, we conducted the ablation experiments for two cases: 1) when the face shape estimation network (FSEN) was removed, 2) when the different face shape estimation network with different keypoint estimation performance was used. For the first case (i.e., Proposed w/o FSEN in Table 6), the experimental settings were set to the same as the previous experimental conditions, and only the existence of the face shape estimation network is different. When the face shape estimation network is removed, the stacked feature map becomes the same as the face appearance feature map (\mathbf{h}_i^F) because the face shape heatmap is not available. According to Table 6, the accuracy dropped by about 2% when the FSEN was removed under the random alignment. Considering that the 1% gain of the recognition performance for the LFW dataset is significant, the FSEN in the proposed method helps improve the recognition performance with robustness. In other words, the explicit face shape prior as well as the well-aligned face image achieves a synergy effect for improving the accuracy to the face misalignment. To deal with the second case, we added a Gaussian random noise for the result of the landmark prediction by the FAN. Let us denote the originally estimated landmark by the FAN as k_x and k_y . Then, the perturbed landmarks by the noise were computed by $\tilde{k}_x = k_x + n_x$ and $\tilde{k}_y = k_y + n_y$, where n_x and n_y follow a Gaussian distribution, i.e., $n_x, n_y \sim \mathcal{N}(0, \epsilon^2)$. In our experiment, we set the standard deviation (ϵ) to 2 and 4 (i.e., Proposed w/ FSEN (FAN $_{\epsilon=2}$) and Proposed w/ FSEN (FAN $_{\epsilon=4}$) in Table 6). Note that the proposed method without the noise perturbation in the FSEN is denoted as Proposed w/ FSEN (FAN $_{\epsilon=0}$). The FSEN performance was measured by the area under the curve (AUC) at 8% normalized mean error (NME) [43] based on the bounding box size normalization [18] with 300W private dataset [44] including 300 indoors and 300 outdoor faces. As shown in Table 6, the accuracy did not change sensitively depending on the face keypoint estimation accuracy. This is mainly because the proposed method utilized the face keypoint heatmap instead of the localization point. Note that the heatmap is for emphasizing the importance of the surrounding areas around the keypoints [32], which can be tolerable to some extent of noise.

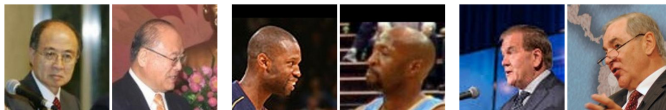
TABLE 7

CPLFW face verification accuracy (%) with or without metadata provided in the dataset. “Difference” (%p) denotes the difference in average accuracy for the two cases.

	CPLFW (1:1 Verification)		
	With metadata	Without metadata	Difference
SphereFace	79.33± 2.45	75.97± 3.42	3.37 ↓
CosFace	85.58± 2.09	79.45± 2.30	6.13 ↓
VGGFace2	83.85± 2.07	79.77± 2.96	4.08 ↓
ArcFace	88.10± 2.58	83.17± 2.53	4.93 ↓
Proposed	85.47± 2.59	85.40± 2.59	0.07 ↓



(a) False Negatives



(b) False Positives

Fig. 8. Example face images for the failure case of the proposed method on the CPLFW dataset.

Experiments with Large Pose Variations. To further investigate the recognition performance under large pose variations where explicit face alignment could not be well addressed, we performed experiments on the CPLFW dataset [45]. As discussed in Section 4.1, the CPLFW dataset is comprised of 3,000 positive and negative pairs (*i.e.*, 6,000 pairs), where the positive pairs with pose differences were collected to add pose variations to intra-class variance. Note that the range of pose variations in the CPLFW is known to be between -90° and $+90^\circ$ in yaw. With the CPLFW dataset, we compared the experimental results for the two cases: 1) “With metadata” that uses the bounding box and landmark information (*i.e.*, metadata) provided by the CPLFW dataset, 2) “Without metadata” that obtains the bounding box and landmark information from the MTCNN detector without the metadata. Table 7 shows the performance and performance differences for each case. We can see that the recognition performance for existing methods decreased by about 4-6%p due to the failure of detecting the bounding box and landmarks. In contrast, the proposed method almost maintains the recognition performance (*i.e.*, -0.07%p difference). This result shows that the proposed method without the explicit face alignment is robust against face misalignment. However, as can be seen in Fig. 8, we can observe several failure cases for the proposed method, where most of the failure cases were those with very large pose variations. This is mainly because the range of pose variations in the VGGFace2 dataset used for training the proposed method is limited (*i.e.*, most of the face images included in VGGFace2 have pose variations between -40° and $+40^\circ$ in yaw).

4.5 Number of Learning Parameters

To validate the efficiency of the proposed method, we measured the number of parameters used for both training and testing stages in Table 8. Among the six different networks listed in Table 1, the parameters for the facial keypoint estimation network occupied 67% of all parameters, which holds the majority number of parameters of the proposed deep network. However, when testing a FR with our method, only two networks (*i.e.*, feature extraction network and feature map embedding network) were used. In other words, 18% parameters (about 6 million) from the total number of training parameters were only used. Therefore, since the proposed deep network only requires a similar number of parameters used in the previous deep FR models without additional computations, we can efficiently compute features robust to the face misalignment.

4.6 Visualization of Activation Maps by Grad-CAM

To understand the effect of the proposed method qualitatively, we visualized the activation maps for input face images. As one of the visualization methods for visual explanations, we adopted the Grad-CAM [47] algorithm that uses the gradients of the target flowing into the final convolutional layer to show the localization map highlighting the important regions [47]. Note that the final convolution layer corresponds to the output of the feature map embedding layer (φ). Specifically, after computing the gradient of the score for the target for the embedded feature maps $\varphi(\mathbf{h}^F)$, the gradients flowing back were globally average-pooled to return the neuron importance weights. Then, the forward activation maps were combined based on the neural importance weights, then activated by the ReLU [47]. The images used for the visualization were obtained from a Web with the same ID label corresponding to the VGGFace2 dataset without face detection, which were not used for training.

Fig. 9 shows the visualization results of the important feature by the Grad-CAM. As shown in the second and third rows, the activations without considering the face shape prior showed lower activation values and were sparsely spread. However, as shown in the fourth and fifth rows, the proposed method trained with the face shape prior concentrated on the important face regions (eyes, nose, and mouth). Less activation occurred in the excluded area that is not closely related face component. Also, the intensity values of the activation map obtained from the proposed method had relatively high values compared to those considering the face feature extraction network only.

4.7 Comparison to the State-of-the-art Methods

Finally, we compared the proposed method to the state-of-the-art methods including both deep FR methods [2], [3], [4], [24], [25], [26], [27], [28], [31], [46] and the face alignment learning algorithms [9], [10], [13], [14] for LFW, YTF, and MegaFace datasets. We brought the FR performance from [2], [24], [25], [26], [27], [28] (for deep FR methods) and [15] (for face alignment learning). In the case of VGGFace2, since the official performance for the MegaFace dataset is not reported in [2], we measured the performance by using the publicly available VGGFace2

TABLE 8

The number of parameters for the proposed method, where the “Train/Test” means that the network parameters are used in both training and test. The ratio is computed based on the number of parameters for the total number of training parameters.

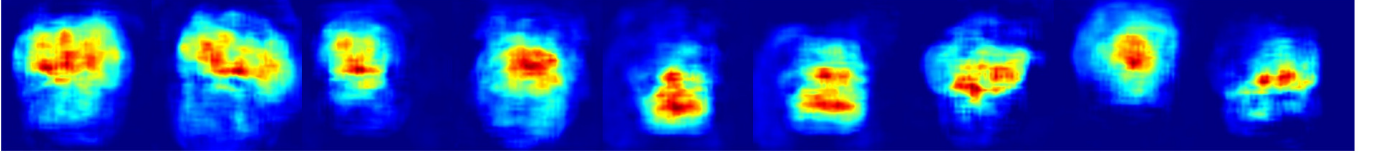
	Train/Test	Number of Parameters	Ratio
Facial Keypoint Estimation Network (S)	Train	23,820,176	67%
Feature Extraction Network (F)	Train/Test	5,902,528	17%
Face Image Decoder (D)	Train	555,712	2%
Fully Connected Layer (Q, b)	Train	4,427,703	12%
Feature Map Embedding Layer (φ)	Train/Test	525,312	1%
Channel Aggregation Layer (ϕ)	Train	526,848	1%
Total Number of Training Parameters	Train	35,758,279	100%
Total Number of Test Parameters	Test	6,427,840	18%



(a) Input images for Grad-CAM visualization



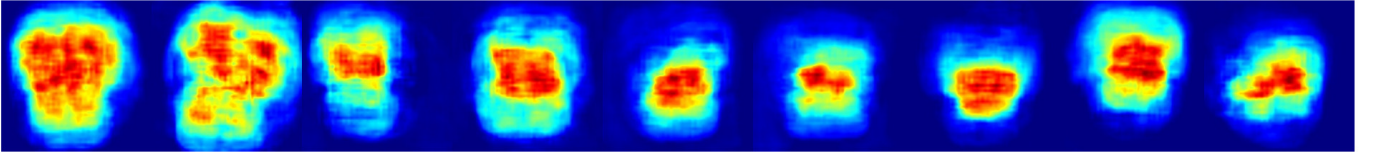
(b) Grad-CAM results without considering the face shape prior (Layered images)



(c) Grad-CAM results without considering the face shape prior (Activations)



(d) Grad-CAM results with considering the face shape prior (Layered images)



(e) Grad-CAM results with considering the face shape prior (Activations)

Fig. 9. Visualization by the Grad-CAM. The top row shows input face images for the analysis. The second and third rows represent the results without considering the face shape prior (face feature extraction network only, *i.e.*, the modified ResNet50). The fourth and fifth rows show the Grad-CAM results considering the face shape prior.

model. For the proposed method, we additionally trained the network based on the ArcFace Loss [2] in order to accommodate the increase of the FR performance, which is denoted as “Proposed + ArcFace Loss”. The residual unit in the face feature extraction network (modified ResNet50) was replaced with the improved residual unit as used in [2].

As can be seen in Table 9, the performance of the deep

FR algorithms shows quite high performance, but its own face alignment algorithm is necessarily required. Since such a specific face alignment algorithm is needed, the performance can be degraded when the alignment type is changed or when the face detection or facial keypoint estimation algorithm fails in a wild environment. In addition, the facial keypoint estimation algorithm for the face alignment

TABLE 9

Comparison with the state-of-the-art deep FR and face alignment learning algorithms, where the FR performance for the MegaFace dataset was measured with the cleaned version of the dataset as suggested in [2]. † denotes the MegaFace performance in [13] without filtering noise lists.

		LFW (1:1 Verification)	YTF (1:1 Verification)	MegaFace (1:N Identification)
Deep Face Recognition	SphereFace [4]	99.42	95.00	72.73 [†]
	CosFace [3]	99.73	97.60	82.72 [†]
	VGGFace2 [31]	99.08	97.30	94.17
	ArcFace [2]	99.82	98.02	98.35
	UniformFace [24]	99.80	97.70	79.98 [†]
	MML [25]	99.63	95.50	83.00 [†]
	DBM [27]	99.78	-	96.35
	CLMLE [26]	99.62	96.50	79.68 [†]
	CurricularFace [28]	99.80	-	98.25
	ARFace [46]	99.62	97.54	96.40
Face Alignment Learning	Zhong <i>et al.</i> [13]	99.33	95.00	65.16 [†]
	ReST [14]	99.03	95.40	-
	GridFace [10]	99.68	95.20	-
	APA [43]	99.68	-	-
	Proposed	99.30	97.00	94.89
	Proposed + ArcFace Loss [2]	99.78	97.90	98.03

requires additional computations. However, the proposed method as a face alignment learning algorithm shows comparable performance without the specific face alignment. Therefore, the proposed method can be used for the FR system efficiently because it is not sensitive to the alignment type and does not require additional computations. In addition, the proposed method outperforms the performance for the face alignment learning methods as shown in Table 9. Therefore, the proposed method would be preferable in terms of effectiveness and efficiency. In other words, the previous face alignment learning algorithms are designed with the face localization network (*e.g.*, recursive spatial transformer modules in ReST [14], rectification network and denoising autoencoder in GridFace [10]) for a face alignment and the feature extraction network for FR, then trained in an end-to-end manner. In the testing phase, the input face image is forwarded to the localization network as well as the feature extraction network. In contrast, since we can extract robust face features to the face misalignment from the face feature extraction network only, our model enables efficient inference in terms of computations and memories.

4.8 Discussion

In this paper, we focus on the observation that the performance of the existing deep face recognition algorithm is degraded if a face image is not well-aligned as used in training. Motivated by recent studies on face alignment learning that learns face alignment and face feature extraction in an end-to-end manner, we propose a face shape-guided face recognition algorithm based on feature alignment (*i.e.*, pixel and feature alignments). Through the experiments with controlled and randomly aligned face images, we observed that the performance of the existing deep face recognition algorithm changed sensitively. In contrast, the proposed method showed robust recognition performance to the face misalignment. This result can be attributed to the proposed method not only being able to see various alignment types during training, but also finding face features by the devised feature alignments using the face shape as a clue (please

refer to Fig. 9). In addition, as shown in Table 9, the proposed method showed comparable performance in a fair experimental environment that each deep face recognition algorithm performs optimal face alignment. Considering that the proposed method is an end-to-end framework integrating both face alignment and face feature extraction, it can be considered efficient. Moreover, we observed that the proposed method outperformed the existing face alignment learning algorithms. In the existing face alignment learning, a localization network in both training and testing is introduced. In the proposed method, the entire network with the help of the localization network (*i.e.*, FAN) is trained via classification and feature alignment losses. However, the localization network is not required in testing. Therefore, our method is computationally efficient even compared to the existing face alignment learning algorithms.

5 CONCLUSION

In this paper, we proposed the face shape-guided deep feature alignment framework for FR robust to the face misalignment. Based on the face shape prior (*i.e.*, face keypoints), we introduced two additional pixel alignment and feature alignment processes with the conventional face feature extraction network, which were learned in an end-to-end manner. For training, the proposed method learned the features for the well-aligned face image by decoding the aggregated features based on a face image and face shape prior. In addition, through the feature alignment process, the learned feature was connected to align with the face shape-guided feature. Through comparative experiments with LFW, CALFW, YTF, and MegaFace datasets, we validated the effectiveness of the proposed method toward face misalignment. In particular, because we do not require additional computations for estimating the face keypoints and the face alignment, it would be efficient for testing a face image.

ACKNOWLEDGMENT

This work was supported by Institute of Information & Communications Technology Planning & Evaluation (IITP)

grant funded by the Korea government (MSIT) (No.2014-3-00123, Development of High Performance Visual BigData Discovery Platform for Large-Scale Realtime Data Analysis and No.2020-0-00004, Development of Previsional Intelligence based on Long-term Visual Memory Network).

REFERENCES

- [1] G. B. Huang, M. Mattar, T. Berg, and E. Learned-Miller, "Labeled Faces in the Wild: A Database for Studying Face Recognition in Unconstrained Environments," in *Technical Report*, 2008.
- [2] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive Angular Margin Loss for Deep Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [3] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "CosFace: Large Margin Cosine Loss for Deep Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [4] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "SphereFace: Deep Hypersphere Embedding for Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2017, pp. 212–220.
- [5] I. Masi, F.-J. Chang, J. Choi, S. Harel, J. Kim, K. Kim, J. Leksut, S. Rawls, Y. Wu, T. Hassner *et al.*, "Learning Pose-Aware Models for Pose-Invariant Face Recognition in the Wild," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 41, no. 2, pp. 379–393, 2018.
- [6] W. Zhang, X. Zhao, J.-M. Morvan, and L. Chen, "Improving Shadow Suppression for Illumination Robust Face Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 41, no. 3, pp. 611–624, 2018.
- [7] Y. Qian, W. Deng, and J. Hu, "Unsupervised face normalization with extreme pose and expression in the wild," in *Proc. of IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 9851–9858.
- [8] H.-I. Kim, S. H. Lee, and Y. M. Ro, "Face Image Assessment Learned with Objective and Relative Face Image Qualities for Improved Face Recognition," in *Proc. IEEE Int'l Conf. Image Processing*, 2015, pp. 4027–4031.
- [9] Z. An, W. Deng, Y. Zhong, Y. Huang, and X. Tao, "APA: Adaptive Pose Alignment for Robust Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition Workshops*, 2019, pp. 227–235.
- [10] E. Zhou, Z. Cao, and J. Sun, "GridFace: Face Rectification via Learning Local Homography Transformations," in *Proc. European Conf. Computer Vision*, 2018, pp. 3–19.
- [11] Y. Wong, C. Sanderson, S. Mau, and B. C. Lovell, "Dynamic Amelioration of Resolution Mismatches for Local Feature based Identity Inference," in *Proc. Int'l Conf. Pattern Recognition*, 2010, pp. 1200–1203.
- [12] J. Y. Choi, Y. M. Ro, and K. N. Plataniotis, "A Comparative Study of Preprocessing Mismatch Effects in Color Image based Face Recognition," *Pattern Recognition*, vol. 44, no. 2, pp. 412–430, 2011.
- [13] Y. Zhong, J. Chen, and B. Huang, "Toward End-to-End Face Recognition through Alignment Learning," *IEEE Signal Processing Letters*, vol. 24, no. 8, pp. 1213–1217, 2017.
- [14] W. Wu, M. Kan, X. Liu, Y. Yang, S. Shan, and X. Chen, "Recursive Spatial Transformer (ReST) for Alignment-Free Face Recognition," in *Proc. IEEE Int'l Conf. Computer Vision*, 2017, pp. 3772–3780.
- [15] H. Wei, P. Lu, and Y. Wei, "Balanced Alignment for Face Recognition: A Joint Learning Approach," *arXiv preprint arXiv:2003.10168*, 2020.
- [16] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A Discriminative Feature Learning Approach for Deep Face Recognition," in *Proc. European Conf. Computer Vision*, 2016, pp. 499–515.
- [17] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint Face Detection and Alignment using Multitask Cascaded Convolutional Networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [18] A. Bulat and G. Tzimiropoulos, "How Far Are We from Solving the 2D & 3D Face Alignment Problem? (and a Dataset of 230,000 3D Facial Landmarks)," in *Proc. IEEE Int'l Conf. Computer Vision*, 2017, pp. 1021–1030.
- [19] A. , Bulat and G. Tzimiropoulos, "Super-FAN: Integrated Facial Landmark Localization and Super-Resolution of Real-World Low Resolution Faces in Arbitrary Poses with GANs," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2018, pp. 109–117.
- [20] M. Zhu, D. Shi, M. Zheng, and M. Sadiq, "Robust Facial Landmark Detection via Occlusion-Adaptive Deep Networks," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 3486–3496.
- [21] H. J. Lee, S. T. Kim, H. Lee, and Y. M. Ro, "Lightweight and Effective Facial Landmark Detection using Adversarial Learning with Face Geometric Map Generative Network," *IEEE Trans. Circuits and Systems for Video Technology*, vol. 30, no. 3, pp. 771–780, 2019.
- [22] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "DeepFace: Closing the Gap to Human-Level Performance in Face Verification," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2014, pp. 1701–1708.
- [23] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2015, pp. 815–823.
- [24] Y. Duan, J. Lu, and J. Zhou, "UniformFace: Learning Deep Equidistributed Representation for Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 3415–3424.
- [25] X. Wei, H. Wang, B. Scotney, and H. Wan, "Minimum margin loss for deep face recognition," *Pattern Recognition*, vol. 97, p. 107012, 2020.
- [26] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Deep Imbalanced Learning for Face Recognition and Attribute Prediction," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 42, no. 11, pp. 2781–2794, 2020.
- [27] D. Cao, X. Zhu, X. Huang, J. Guo, and Z. Lei, "Domain Balancing: Face Recognition on Long-Tailed Domains," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2020, pp. 5671–5679.
- [28] Y. Huang, Y. Wang, Y. Tai, X. Liu, P. Shen, S. Li, J. Li, and F. Huang, "CurricularFace: Adaptive Curriculum Learning Loss for Deep Face Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2020, pp. 5901–5910.
- [29] A. Newell, K. Yang, and J. Deng, "Stacked Hourglass Networks for Human Pose Estimation," in *Proc. European Conf. Computer Vision*, 2016, pp. 483–499.
- [30] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial Transformer Networks," in *Proc. Advances in Neural Information Processing Systems*, 2015, pp. 2017–2025.
- [31] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "VG-GFace2: A Dataset for Recognising Faces across Pose and Age," in *Proc. IEEE Int'l Conf. Automatic Face and Gesture Recognition*, 2018, pp. 67–74.
- [32] P. Khorramshahi, A. Kumar, N. Peri, S. S. Rambhatla, J.-C. Chen, and R. Chellappa, "A Dual-Path Model With Adaptive Attention for Vehicle Re-Identification," in *Proc. IEEE Int'l Conf. Computer Vision*, 2019, pp. 6132–6141.
- [33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [34] Y. Suh, J. Wang, S. Tang, T. Mei, and K. Mu Lee, "Part-aligned Bilinear Representations for Person Re-Identification," in *Proc. European Conference on Computer Vision*, 2018, pp. 402–419.
- [35] B. Bozorgtabar, M. S. Rad, D. Mahapatra, and J.-P. Thiran, "SynDeMo: Synergistic Deep Feature Alignment for Joint Learning of Depth and Ego-Motion," in *Proc. IEEE Int'l Conf. Computer Vision*, 2019, pp. 4210–4219.
- [36] M. Lin, Q. Chen, and S. Yan, "Network in Network," in *Proc. Int'l Conf. Learning Representations*, 2013.
- [37] T. Zheng, W. Deng, and J. Hu, "Cross-Age LFW: A Database for Studying Cross-Age Face Recognition in Unconstrained Environments," *arXiv preprint arXiv:1708.08197*, 2017.
- [38] L. Wolf, T. Hassner, and I. Maoz, "Face Recognition in Unconstrained Videos with Matched Background Similarity," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011, pp. 529–534.
- [39] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard, "The MegaFace Benchmark: 1 Million Faces for Recognition at Scale," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 4873–4882.
- [40] H.-W. Ng and S. Winkler, "A Data-Driven Approach to Cleaning Large Face Datasets," in *Proc. IEEE Int'l Conf. Image Processing*, 2014, pp. 343–347.
- [41] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning Face Representation from Scratch," *arXiv preprint arXiv:1411.7923*, 2014.
- [42] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "MS-Celeb-1M: A Dataset and Benchmark for Large-Scale Face Recognition," in *Proc. European Conf. Computer Vision*, 2016, pp. 87–102.

- [43] X. Wang, L. Bo, and L. Fuxin, "Adaptive Wing Loss for Robust Face Alignment via Heatmap Regression," in *Proc. IEEE Int'l Conf. Computer Vision*, 2019, pp. 6971–6981.
- [44] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic, "300 Faces in-the-Wild challenge: The First Facial Landmark Localization Challenge," in *Proc. IEEE Int'l Conf. Computer Vision Workshops*, 2013, pp. 397–403.
- [45] T. Zheng and W. Deng, "Cross-Pose LFW: A Database for Studying Cross-Pose Face Recognition in Unconstrained Environments," *Beijing University of Posts and Telecommunications, Tech. Rep*, vol. 5, p. 7, 2018.
- [46] L. Zhang, L. Sun, L. Yu, X. Dong, J. Chen, W. Cai, C. Wang, and X. Ning, "ARFace: Attention-aware and Regularization for Face Recognition with Reinforcement Learning," *IEEE Trans. Biometrics, Behavior, and Identity Science*, 2021.
- [47] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization," in *Proc. IEEE Int'l Conf. Computer Vision*, 2017, pp. 618–626.