

# Information Maximization for Extreme Pose Face Recognition

Mohammad Saeed Ebrahimi Saadabadi, Sahar Rahimi Malakshan,  
Sobhan Soleymani, Moktari Mostofa, and Nasser M. Nasrabadi

me00018, sr00033, ssoleyma, mm0251@mix.wvu.edu, nasser.nasrabadi@mail.wvu.edu

## Abstract

In this paper, we seek to draw connections between the frontal and profile face images in an abstract embedding space. We exploit this connection using a coupled-encoder network to project frontal/profile face images into a common latent embedding space. The proposed model forces the similarity of representations in the embedding space by maximizing the mutual information between two views of the face. The proposed coupled-encoder benefits from three contributions for matching faces with extreme pose disparities. First, we leverage our pose-aware contrastive learning to maximize the mutual information between frontal and profile representations of identities. Second, a memory buffer, which consists of latent representations accumulated over past iterations, is integrated into the model so it can refer to relatively much more instances than the mini-batch size. Third, a novel pose-aware adversarial domain adaptation method forces the model to learn an asymmetric mapping from profile to frontal representation. In our framework, the coupled-encoder learns to enlarge the margin between the distribution of genuine and imposter faces, which results in high mutual information between different views of the same identity. The effectiveness of the proposed model is investigated through extensive experiments, evaluations, and ablation studies on four benchmark datasets, and comparison with the compelling state-of-the-art algorithms.

## 1. Introduction

With the advancement of technology and increasing demand for security, biometrics are among the most essential and surfed applications of computer vision [25]. Among biometric traits, the face has received particular attention since it is naturally exposed, offers better hygiene in the acquisition, and can be acquired in an unconstrained setting without direct participation of the user [23]. Face Recognition (FR) has been a major interest in computer vision for many years, and FR methods have advanced significantly over the years [23]. Classical FR techniques are mainly

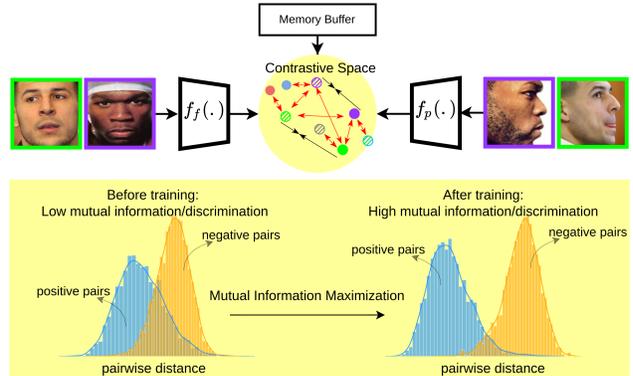


Figure 1. (Top) Two pairs of positive samples are fed to the coupled-encoder ( $f_f(\cdot)$  and  $f_p(\cdot)$ ). Each color in the contrastive space presents a distinct identity. Black arrows demonstrate the attraction between two representations, and the red arrows increase the distance. Solid and dashed circles represent profile and frontal representations, respectively. Due to the memory buffer, the number of instances in the contrastive space is more than the mini-batch size. (Bottom) Illustrating the distance distributions between positive (shown in blue) and negative (shown in orange) pairs. Our optimization improves the similarity between the different views of the same identity while increasing the distance between different identities.

based on extracting hand-crafted features, and the primary concern is extracting features with high intra-class compactness, and inter-class separability [2].

Modern approaches address this issue by incorporating learning models based on the Convolutional Neural Networks (CNNs) [35]. CNNs have demonstrated extraordinary performance in FR; however, their performance drastically degrades for profile views [56]. There are four primary issues with profile face images as compared to frontal images: 1) self-occlusion, 2) background distraction, 3) shift in the distribution of data, and 4) inaccuracy in the alignment due to the lack of accurate landmarks [4, 38]. There are two mainstream approaches for profile-to-frontal FR. The first approach handles pose variation by extracting pose-invariant features [21, 45]. The second approach estimates the frontal view of a given profile face and then

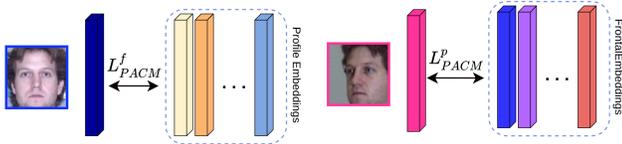


Figure 2. With memory, for each sample in mini-batch  $B$ , we can have different contrastive pairs and calculate the loss with the samples which are not in the current mini-batch. Different colors refer to representation of different identity in the dataset

utilizes it for the recognition task [40, 15].

In the first approach, Softmax with a cross-entropy loss is mainly adopted to supervise a deep classifier. Although Softmax promotes the separability of representations, provided features are insufficiently discriminative in practical FR problems [45]. To address this issue, pioneering works of [35, 48, 37] employ sample-to-sample comparison as their loss functions to reduce the intra-class variations. However, most of the recent FR methods mainly focus on the sample-to-prototype comparison, and they improve the discriminative power of representations by applying several margin penalties on the Softmax loss function [21, 45]. Despite the remarkable performance, a common issue of these approaches is that FR datasets contain a large number of identities, and only a few identities are presented in each mini-batch, which complicates finding an optimal decision boundary [12]. Increasing the mini-batch size may alleviate the problem. However, it does not guarantee performance improvement [16, 53], and it may not be possible (due to the memory constraints). Moreover, shortcomings such as sensitivity to noisy labels [55], likelihood of poor margin [9], and convergence difficulty on the networks with small embedding feature size [18] have led to diminishing generalization.

In the second approach, FR process is separated into two tasks: identity-preserving face generation and frontal face recognition. Among face generation modules, the Generative Adversarial Networks (GANs) have received special attention [40, 15]. Despite the remarkable results concerning image quality and human perception, GANs add high-frequency components to the synthesized images, which negatively affects the recognition process [46]. Besides, from the optimization perspective, profile-to-frontal face generation is an intrinsically ill-posed problem, and multiple frontal faces exist for each profile face [14]. Also, there are several other nuisance factors in face images, including expression, illumination, and the quality of images. These factors result in a large gap between features of real and synthesized frontal faces in the identity metric space, which significantly deteriorates the final performance [46].

In this paper, we hypothesize that profile and frontal faces have latent connection in an abstract embedding space. We exploit this hidden connection in the embedding

domain using a deep coupled model consisting of dedicated networks for the profile and frontal views of the face. These two networks share the same discriminative latent embedding, see Fig. 1. Using the proposed Pose-Aware Contrastive learning (PAC), we enforce the agreement of the features by maximizing the lower bound of the mutual information between the representations of the same identity [39]. In this manner, the model aims to pull closer the representations from pairs of the same identity compared to representations of different identities [17]. PAC also helps the model to implicitly benefit from hard negative/positive instances [17]. Hard samples are close to the decision boundary in the embedding space, and emphasizing them in training leads to faster convergence, and better generalization [35]. Also, we leverage Pose-Aware Contrastive with Memory buffer (PACM), a simple yet effective way to help the loss utilize a massive number of identities’ representations without increasing the mini-batch size.

Aiming to further reduce the gap between the profile and frontal images in the embedding space, we employ our proposed Pose-Aware Adversarial Domain Adaptation (PADA) learning approach to enforce the model to learn an asymmetric mapping from profile to frontal representation. Our experiments, evaluations, and ablations studies show that the proposed framework achieves notable performance in learning pose-invariant discriminative representations. Contributions of this paper can be summarized as follows:

- A novel profile-to-frontal face recognition model is developed, which utilizes a pose-aware contrastive learning to maximize the mutual information between the profile and frontal representations from the same identity in an embedded space.
- A novel pose-aware domain adaptation approach is developed to enforce the agreement of features from different poses.
- A novel approach is proposed to learn pose-agnostic representations from a larger number of instances than the mini-batch size in a multiview setting.

## 2. Related Works

Deep learning has been applied to various applications since its advent [7, 29, 30, 1, 34, 27, 28]. Biometric has been one the most surfed area due to the availability of large-scale datasets which can be either in the form of signals or images. Among biometric traits, the face has received special attention. In this section, we briefly summarize recent attempts in FR.

## 2.1. Deep Face Recognition

The availability of computing power has made CNN the primary tool in computer vision [23]. During the past decade, the introduction of new network architectures, accessibility to large-scale datasets, and modifications of the loss functions have led to significant achievements in deep FR [23, 5]. Along with other supervised deep learning frameworks, Softmax with a cross-entropy loss is of the most popular criterion for FR [17]. Intrinsically, features provided by the Softmax have angular distribution [7, 48]. Studies have shown that considering angular distance instead of Euclidean distance significantly improves the FR performance [45, 48]. Based on this characteristics, multiple training paradigms have been proposed to adapt various kinds of margins to the Cosine based embedding space [45, 48, 21]. Sample-to-sample loss functions are also well established in deep FR [37, 48, 35]. Sun *et al.* [37] combined the identification and Margin Contrastive Loss (MCL) loss for having a more powerful supervisory signal. In [35], Schroff *et al.* presented the Triplet loss, which forces the representations of a triplet to be discriminative. To increase intra-class compactness, Wen *et al.* [48] introduced the Center loss in which the model learns a center for each class. They used it in combination with the Softmax loss to keep representations from collapsing to zero.

## 2.2. Pose Robust Face Recognition

Although near-frontal FR is considered a solved problem in common cases, FR in extreme poses, where enrolled faces in the gallery and the probe images have large pose disparity, has still remained a difficult effort [42]. There are two main approaches to cope with profile-to-frontal FR [14, 13, 22, 26]. One major research avenue is based on synthesizing a frontal face from its profile input and then utilizing the synthesized frontal face for recognition [14, 13, 32, 52, 40]. Despite the satisfactory results, this approach has a handful of intrinsic drawbacks. First, the problem of recovering a face’s canonical view from its profile pair is under-defined [56, 14]. Second, since the frontal faces should be generated first in order to train the classifier, end-to-end training of the generator and classifier is unattainable [42].

Another main line of inquiry for profile-to-frontal FR is to learn pose-agnostic mapping. For instance, Masi *et al.* [22] used 3D rendering to synthesize multiple views of a profile image. These images were used to train pose-specific deep feature extractors. Then, features from all networks are fused to construct the final prediction. DREAM [4] is based on finding a mapping between profile and frontal embeddings. The authors hypothesized a gradual connection from profile to frontal representation. They utilized a residual mapping that adds pose-adaptive residuals to the features extracted from profile faces. Meng *et al.*

[26] disentangled pose and identity representations by mapping face to identity and landmark subspaces. To disentangle identity from the pose, Yin *et al.* adopted a multi-tasking framework [51]. PF-cpGAN [38] seeks to learn pose-agnostic representation by employing face frontalization as a sub-task.

Although these methods have shown promising results, there is still a significant performance gap for faces with extreme poses in unconstrained conditions [56]. In comparison, our proposed coupled-encoder benefits from PACM and PADA losses. They encourage coupled-encoder to map the faces with the same identity to close representations and faces with different identities to representations far from each other. In addition, the memory buffer elevates the efficiency of proposed contrastive learning by looking at a relatively much larger number of identities than the mini-batch size.

## 3. Proposed Method

We introduce a coupled-encoder architecture to map the profile and frontal face images to a shared embedding space. PACM and PADA losses force the coupled-encoder to learn pose-agnostic representations. Our model incorporates a massive number of instances to PACM loss function. For example, on a single NVIDIA TITAN X GPU and the mini-batch size of 32, the coupled-encoder calculates the loss between more than 6000 distinct instances, which is beneficial for contrastive learning to maximize the mutual information between two different views of the face images of an identity [44]. Furthermore, PADA loss improves the compatibility of representations by forcing the profile encoder to map the off-angle faces close to its frontal pairs.

### 3.1. Pose-Aware Contrastive Learning

Our proposed profile-to-frontal FR framework is based on recent information-theoretic techniques for contrastive representation learning [44, 39]. During training, we aim to maximize the mutual information between profile and frontal face images from the same identity (positive pair) and minimize the mutual information between face images with different identities (negative pairs). During testing, we make the decision considering the distance between representations of a pair of face images in the embedding space. The first step during training is to select face images that represent the same identity but distinct views. Then, each image in positive pair is employed to choose instances that represent distinct identities and views. For example, given a profile face image, we should pick negative frontal instances and vice versa. Then, these pairs of frontal and profile face images are used to train two deep encoders.

Given a training set  $D = D_f \cup D_p$ ,  $D_f : \{(x_{f,i}, y_{f,i})\}_{i=0}^{N_f}$  and  $D_p : \{(x_{p,i}, y_{p,i})\}_{i=0}^{N_p}$  represent profile and frontal subsets of dataset, respectively.

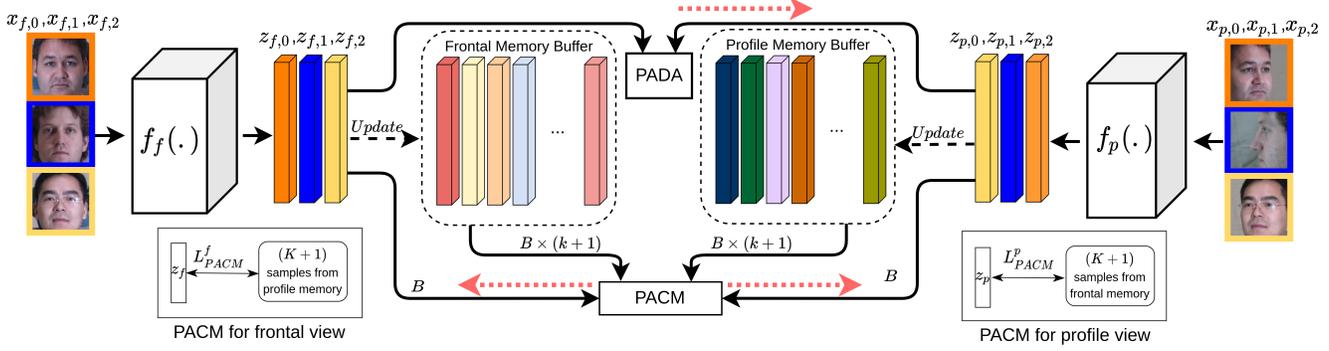


Figure 3. A mini-batch ( $B$ ) of frontal and profile faces is fed into the coupled-encoder. This model provides representations with high mutual information (high similarity in the embedding space) for genuine pairs and far distant representations for the imposter pairs. Here, positive pairs are shown in the same color. Without memory buffer, the contrastive loss is calculated within the mini-batch (two negative samples for each image). However, memory buffer provide  $K$  negative samples,  $K > B$ , for each sample. Red-dashed arrows show the gradient back-propagation. The pose-aware adversarial domain adaptation (PADA) loss only affects the profile encoder.

$N_p$  is the number of profile samples and  $N_f$  frontal samples. As presented in Fig. 3, the coupled-encoder consists of  $f_f(\cdot)$  and  $f_p(\cdot)$ , which are the frontal and profile dedicated embedding sub-networks. These two sub-networks map the frontal and profile faces to a  $d$ -dimensional embedding space:  $f_f(\cdot) : \mathbb{R}^{3 \times w \times h} \rightarrow \mathbb{R}^d$  and  $f_p(\cdot) : \mathbb{R}^{3 \times w \times h} \rightarrow \mathbb{R}^d$ , respectively. For convenience, we ignore the index  $i$  that reflects the index of the sample in their corresponding subset.  $z_f = f_f(x_f)$  and  $z_p = f_p(x_p)$  are the representations of the frontal,  $x_f$ , and profile,  $x_p$ , images generated by their corresponding encoders.

The first step toward our training paradigm is to select pair of images with the same identity and pose disparity. To construct our training samples, we define a genuine (positive) pairs as:

$$P = \{(z_f, z_p) | y_f = y_p\}, \quad (1)$$

where  $y_f$  and  $y_p$  represent labeled identities of  $z_f$  and  $z_p$ , respectively. In contrast, we define an imposter (negative) pair as:

$$N = \{(z_f, z_p) | y_f \neq y_p\}, \quad (2)$$

for positive pairs, we first choose a random identity  $y_0$  as an anchor. Then, sampling two images independently from  $D_f$  and  $D_p$  given the selected identity. Consequently, we define the joint distribution of genuine pair as [44]:

$$\begin{aligned} p(z_f, z_p) &= \sum_{y_0 \in C} p(y_f = y_0) p(z_f | y_f = y_0) \\ &\quad \times p(y_p = y_0) p(z_p | y_p = y_0) \\ &= \sum_{y_0 \in C} p(y_f = y_0, y_p = y_0) \\ &\quad \times p(z_f, z_p | y_f = y_p = y_0) \\ &= \sum_{y_0 \in C} p(z_f, z_p, y_f = y_0, y_p = y_0), \end{aligned} \quad (3)$$

where  $C$  reflects the total identities presented in the dataset. Assuming a high entropy of negative pairs and the large number of identities, which is the case for the FR, we approximate the sampling of the negative pairs with sampling from product of marginals [44]:

$$\begin{aligned} p(z_f) p(z_p) &\approx \sum_{y_0 \in C} \sum_{\substack{y'_0 \in C \\ y_0 \neq y'_0}} p(y_f = y_0) p(z_f | y_f = y_0) \\ &\quad \times p(y_p = y'_0) p(z_p | y_p = y'_0). \end{aligned} \quad (4)$$

We aim to train  $f_f(\cdot)$  and  $f_p(\cdot)$  such that face images of an identity in different views are mapped closely in the embedding space. To this end, we maximize the mutual information between positive pair representations by maximizing the KL divergence between Eqs. 3 and 4 [39]. Hence, we aim to learn a function,  $h(\cdot)$ , which provides a low value for negative pairs and high value for positive pairs [44].

$$h(z_f, z_p) = \exp\left(\frac{1}{\tau} \frac{z_f \cdot z_p}{\|z_f\| \cdot \|z_p\|}\right), \quad (5)$$

$h(\cdot)$  reflects the cosine similarity between latent representations and  $\tau$  is the temperature [11], which plays an important role in concentration of representations in the hypersphere [50, 45].

The contrastive learning aims to pull an anchor  $z_{f,i_0}$  and positive samples  $z_{p,i_0}$  close in the embedding space while pushing the anchor away from many negative samples [39]:

$$L_{cont} = -\mathbb{E}_S \left[ \log \frac{h(z_{f,i_0}, z_{p,i_0})}{\sum_{i=0}^k h(z_{f,i_0}, z_{p,i})} \right], \quad (6)$$

where  $S : \{(z_{f,i_0}, z_{p,i})\}_{i=0}^k$  is a set of  $k$  negative pairs and one positive pair [44]. In [44], it is proven that optimal  $h(\cdot)$  is in direct proportion with the ratio of joint distribution and

Table 1. Verification accuracy (%) and standard deviation for CFP-FP over standard 10-folds. Results of the [4, 6] are copied from [38].

Method	Frontal-Profile		Frontal-Frontal	
	Accuracy	EER	Accuracy	EER
PR-REM [4]	93.25(2.23)	7.92(0.98)	98.1(2.19)	1.1(0.22)
DCNN [6]	84.91(1.82)	14.97(1.98)	96.4(0.69)	3.48(0.67)
p-CNN [51]	94.39(1.17)	5.94(0.11)	97.79(0.40)	2.48(0.07)
FRN-TI [42]	95.62	-	-	-
DR-GAN [40]	93.41(1.17)	-	97.84(0.79)	-
PF-cpGAN [38]	93.78(2.46)	7.21(0.65)	98.88(1.56)	0.93(0.14)
PIM [56]	93.1(1.01)	7.69(1.29)	<b>99.44(0.36)</b>	0.86(0.49)
ours	<b>95.85(1.07)</b>	<b>4.22(0.15)</b>	99.37(0.4)	<b>0.63(0.05)</b>

product of marginals distributions:  $h \propto \frac{p(z_f, z_p)}{p(z_f)p(z_p)}$ . Replacing  $h(\cdot)$  with the density ratio results [44]:

$$L_{cont}^{optim} \geq \log(k) - \mathbb{E}_{(z_f, z_p) \sim p_{z_f, z_p}} \log \left[ \frac{p(z_f, z_p)}{p(z_f)p(z_p)} \right], \quad (7)$$

recalling the mutual information between two random variable  $z_f$  and  $z_p$ :  $I(z_f; z_p) = \mathbb{E}_{p_{z_f, z_p}} \left[ \frac{p(z_f, z_p)}{p(z_f)p(z_p)} \right]$ . Consequently, for any positive pair of frontal and profile faces:  $I(z_f, z_p) \geq \log(k) - l_{cont}^{optim}$  [44]. Hence, minimizing the contrastive loss results in maximizing the lower bound to  $I(z_f; z_p)$ .

Without loss of generality, we consider the features are normalized:  $\|z_f\| = \|z_p\| = 1$ . Consequently, given a mini-batch, the PAC loss for a frontal anchors is:

$$L_{PAC}^f = - \sum_{i=1}^{|B|} \log \frac{\exp(\frac{1}{\tau} z_{f,i} \cdot z_{p,a_i})}{\sum_{j \in N_p(i)} \exp(\frac{1}{\tau} z_{f,i} \cdot z_{p,j})}, \quad (8)$$

where  $B$  is the mini-batch and  $N_p(i)$  is a set of one positive and many negative profile samples corresponding to  $z_{f,i}$ . Symmetrically, considering profile samples as anchors:

$$L_{PAC}^p = - \sum_{i=1}^{|B|} \log \frac{\exp(\frac{1}{\tau} z_{p,i} \cdot z_{f,a_i})}{\sum_{j \in N_f(i)} \exp(\frac{1}{\tau} z_{p,i} \cdot z_{f,j})}, \quad (9)$$

there are four main advantages in using this loss function. 1) Comparison with every negative sample within mini-batch at the same time, denominator in Eqs. 8 and 9, 2) rather than optimizing the angle between the representations and their corresponding prototypes, the model directly learns the angle between representations [21], 3) PAC provides the implicit hard negative/positive mining to the model [17], and 4) PAC maximizes the mutual information between representations from different views of the shared context [3].

### 3.2. Pose-Aware Contrastive Learning with Memory Buffer

Due to the large number of classes in FR datasets and having a small number of identities within a mini-batch,

the conventional FR methods could not cover a large number of identities at each step of loss calculation [12]. Increasing the mini-batch size may alleviate the issue, but it does not ensure the improvement in the performance, and in many cases, due to the memory constraint, it is impractical [16, 53]. We mitigate this issue by adopting the memory buffer framework [50] to the loss function to benefit from more negative instances. The memory buffer consists of the latent representations of profile and frontal faces from past iterations. During each learning iteration, representations  $z_f$  and  $z_p$  are updated to the memory at the corresponding instance entry:

$$\begin{aligned} r_{f,t+1} &= m * r_{f,t} + (1 - m) * z_f \\ r_{p,t+1} &= m * r_{p,t} + (1 - m) * z_p, \end{aligned} \quad (10)$$

where  $r_f$  and  $r_p$  stand for features saved in frontal and profile memory buffer and  $m$  is the momentum coefficient in updating [50]. Therefore, Eq. 8 can be rewritten as:

$$L_{PACM}^f = - \sum_{i=1}^{|B|} \log \frac{\exp(\frac{1}{\tau} z_{f,i} \cdot r_{p,a_i})}{\sum_{j \in N_p^M(i)} \exp(\frac{1}{\tau} z_{f,i} \cdot r_{p,j})}, \quad (11)$$

where  $N_p^M(i)$  represent a set of one positive and multiple negative profile representations drawn from the memory. Similarly:

$$L_{PACM}^p = - \sum_{i=1}^{|B|} \log \frac{\exp(\frac{1}{\tau} z_{p,i} \cdot r_{f,a_i})}{\sum_{j \in N_f^M(i)} \exp(\frac{1}{\tau} z_{p,i} \cdot r_{f,j})}. \quad (12)$$

Instances within a mini-batch form the  $N_f(i)$  and  $N_p(i)$ ; however,  $N_f^M(i)$  and  $N_p^M(i)$  are drawn from the memory and are not constrained to the mini-batch size. Therefore, we can chose the number of negative pairs such that:  $|N_p^M(i)| \gg |N_p(i)|$  and  $|N_f^M(i)| \gg |N_f(i)|$ . Memory allows us to choose different samples for every instance in each mini-batch and not be limited to the mini-batch size [41], see Fig. 2. Finally, the overall loss for the Pose-Aware Contrastive with Memory buffer (PACM) is:

$$L_{PACM} = L_{PACM}^f + L_{PACM}^p. \quad (13)$$

### 3.3. Pose-Aware Adversarial Domain Adaptation Learning

For each identity, the ideal scenario is to have an identical profile and frontal feature representations. To further improve the similarity of these features, we adapt the idea of adversarial adaptation [43], which aims at making the representations of profile and frontal images as similar as possible. To this end, we aim to fool a binary classifier (view discriminator), which is going to be trained to distinguish between profile and frontal representations. This mimics the policy which is used in the GANs in which the generator tries to produce samples that are indistinguishable

Table 2. Performance (%) comparison of our framework on Setting1 and Setting2 of Multi-PIE dataset. Last three rows represent our results with single (first two) and couple network (last row) with 200 negative samples.

Method	Setting 1						Setting 2					
	$\pm 15^\circ$	$\pm 30^\circ$	$\pm 45^\circ$	$\pm 60^\circ$	$\pm 75^\circ$	$\pm 90^\circ$	$\pm 15^\circ$	$\pm 30^\circ$	$\pm 45^\circ$	$\pm 60^\circ$	$\pm 75^\circ$	$\pm 90^\circ$
FF-GAN [52]	-	-	-	-	-	-	94.6	92.5	89.7	85.2	77.2	61.2
DR-GAN [40]	-	-	-	-	-	-	94.0	90.1	86.2	83.2	-	-
FNM+Light CNN [32]	99.9	99.5	98.2	93.7	81.3	55.8	-	-	-	-	-	-
CAPG-GAN [13]	99.95	99.37	98.28	93.74	87.40	77.10	99.82	99.56	97.33	90.63	83.05	66.05
TP-GAN [14]	99.78	99.85	98.58	92.93	84.10	64.03	98.68	98.06	95.38	87.72	77.43	64.64
PIM [56]	99.80	99.40	98.30	97.70	91.20	75.00	99.30	99.00	98.50	98.10	95.00	86.50
PoseFace [26]	<b>100</b>	99.97	99.62	98.55	96.07	90.58	-	-	-	-	-	-
FFWM [47]	<b>100</b>	<b>100</b>	<b>100</b>	98.86	96.54	88.55	99.86	99.80	99.37	98.85	97.20	93.17
PF-cpGAN [38]	99.9	99.9	98.9	97.6	94.2	88.1	-	-	-	-	-	-
Ours-single-wo PADA	<b>100</b>	<b>100</b>	99.97	99.53	97.26	91.70	99.98	99.97	99.82	99.47	97.19	91.51
Ours-single	<b>100</b>	<b>100</b>	99.97	99.63	97.50	92.65	99.95	99.83	99.70	98.97	96.64	91.79
Ours-couple	<b>100</b>	<b>100</b>	99.97	<b>99.74</b>	<b>98.04</b>	<b>94.64</b>	<b>100</b>	<b>100</b>	<b>99.98</b>	<b>99.76</b>	<b>98.21</b>	<b>94.49</b>

from real samples [43]. The view discriminator  $D$  classifies whether a feature vector comes from a profile or frontal image. Therefore,  $D$  should maximize cross-entropy loss function:

$$L_D(x_p, x_f, f_p(\cdot), f_f(\cdot)) = \mathbb{E}[\log D(f_f(x_f))] + \mathbb{E}[\log(1 - D(f_p(x_p)))]. \quad (14)$$

It is important to remember that for adversarial domain adaptation learning, positive samples are used,  $y_f = y_p$ . The profile encoder’s objective is to fool the view discriminator by maximizing the following:

$$L_{encoder}(x_p, f_p(\cdot)) = \mathbb{E}[\log D(f_p(x_p))], \quad (15)$$

considering Eqs. 14 and 15, we conclude that the profile encoder and the view discriminator play a minmax game:

$$L_{PADA} = \min_{f_p(\cdot)} \max_D \{ \mathbb{E}[\log D(f_f(x_f))] + \mathbb{E}[\log(1 - D(f_p(x_p)))] \}, \quad (16)$$

where  $f_f(\cdot)$  is fixed during the adversarial training and  $f_p(\cdot)$  is trained by Eq. 16. Consequently, the model learns an asymmetric mapping from profile to frontal representation that modifies the profile encoder to learn representations that match the representations of frontal faces [43].

Based on the above discussion, we formulate the total training loss function as:

$$L_{total} = \lambda_1 l_{PADA} + \lambda_2 l_{PACM}, \quad (17)$$

where  $\lambda_1$  and  $\lambda_2$  are the training regularization parameters.

## 4. Experiments

We study the performance of the coupled-encoder on four FR datasets. We report the results of the proposed framework for verification and identification setup and compare them with the state-of-the-art (SOTA) methods.

Furthermore, we investigate the impact of the number of negative samples and effect of different terms in Eq. 17.

### 4.1. Training Setup

For all datasets, MTCNN [54] is considered to detect and align faces. All the images are resized to  $112 \times 112$ , and pixel values are normalized to  $[-1, 1]$ . Most of the FR datasets are imbalanced in two aspects: 1) the number of per identity samples, and 2) the number of profile and frontal samples for each identity. Consequently, there is a good chance that one of the networks learns a degenerate solution [8]. We initialize the encoders sub-networks with pre-trained weights on the VggFace2 dataset [5] with the Softmax loss to mitigate this problem.

For selecting the frontal and profile pairs, we apply [33] to the datasets to create frontal and profile subsets based on the yaw angle. Then, face images with an absolute yaw value less than  $15^\circ$  are considered frontal. For training Eq. 17, the initial learning rate is set to 0.001 and is multiplied by 0.1 every ten epochs, and weight decay and momentum are 0.00001 and 0.9, respectively. The model is trained for 20 epochs. During the training,  $\lambda_1 = 0.1$  and  $\lambda_2 = 1.0$ , and the number of negative samples is 6,000. We adopt ResNet50 [7] as the encoder networks for the profile and frontal views. The final feature is of size  $512 \times 7 \times 7$ . Feature maps are reshaped to form a vector of size 25,088 and passed to a fully-connected layer of size 512 to construct the final representation. The last feature vectors of encoders are enqueued to the frontal and profile memory buffer, see Fig. 3. The view discriminator is an MLP with two hidden layers of size 256, each followed by batch normalization and leaky relu activation function. At the top of these hidden layers, there is a single neuron with a sigmoid activation function. The model is trained using Stochastic Gradient Descent (SGD) with a mini-batch size of 32 on an

## 4.2. Results

**CMU Multi-PIE** dataset [10] includes 750,000 images of 337 identities with variations in pose, illumination, and expression in the controlled environment. It contains images of 15 different views from 20 illuminations in different expressions. For fair comparison, we use neutral images in 13 views of  $\{0^\circ, \pm 15^\circ, \pm 30^\circ, \pm 45^\circ, \pm 60^\circ, \pm 75^\circ, \pm 90^\circ\}$ , and all the variations in illumination are included as [14]. More specifically, there are two main settings for this dataset in the literature: Setting1 and Setting2. In both settings, face images with neutral expression are used. In the Setting1, 250 identities from session 01 of the dataset are employed. The first 150 identities are selected for training and the rest for testing. The test set consists of probe and gallery sets. The gallery set consists of one frontal image per identity in neutral illumination, and the probe set contains images with a yaw angle other than zero. There is no overlap between training and testing identities. In Setting2, face images of the 200 identities from four sessions are used for training. The probe and gallery sets for testing are constructed as Setting1. In this dataset, images with zero yaw degree are considered frontal, and all the other views are considered profile. Following [40, 32, 14, 56, 26, 47, 38], we fine-tune the coupled-encoder on the training sets of Setting1 and Setting2, separately.

Table 2 demonstrates the performance of our model in comparison with SOTA models on the Multi-PIE dataset and we investigate the effect of face angle on the FR performance. Almost all methods perform above 92% for pose variation in  $[-60^\circ, +60^\circ]$ . However, beyond this range, their performance drastically decreases. Table 2 shows that the coupled-encoder outperforms both face frontalization and pose-invariant feature learning algorithms for the face with  $\pm 90^\circ$  pose. Our framework improves [38] by almost 6%. Even without the pose-aware adversarial domain adaptation, the presented model achieves better results compared to other methods (in Setting1). Considering PoseFace [26], coupled-encoder performs better for every pose, and it outperforms PoseFace for face images with  $\pm 90^\circ$  pose by almost four percentage points. Same as Setting1, the coupled-encoder surpasses other algorithms in Setting2. The improvement is more noticeable for the extreme pose of  $\pm 90^\circ$ .

The **Celebrities in Frontal-Profile in the Wild (CFP)** dataset [36] includes facial images of 500 different identities with 10 frontal and 4 profile samples per identity. Following [36], we evaluate our method on 10-fold protocol, and each fold consists of 350 genuine and 350 imposter pairs. Considering results on the CFP-FP dataset in Table 1, our framework performs better than the other SOTA methods with at least a margin of 0.23% in terms of accuracy and 1.72% improvements in terms of EER.

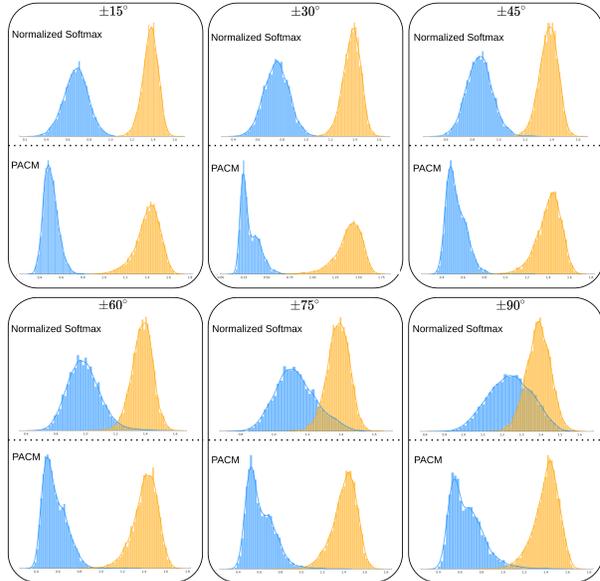


Figure 4. The distance distributions between positive (shown in blue) and negative (shown in orange) pairs. In near frontal poses, the normalized Softmax (top in each block) and the PACM (bottom in each block) separate positive from negative pairs. In near-profile,  $\pm 60^\circ$  and  $\pm 75^\circ$ , and complete profile,  $\pm 90^\circ$ , only PACM perfectly separates these two distributions (High mutual information between positive pairs).

Table 3. Verification accuracy (%) for IJB-B and IJB-C dataset. Results of the [35, 32, 4] are copied from [38]

Method	IJB-C(TAR@FAR)		IJB-B(TAR@FAR)	
	0.001	0.01	0.001	0.01
GOTs [49, 24]	36.3	62.1	33.0	60.0
VGG-CNN [49, 24]	74.3	87.2	72.0	86.0
CFR-GAN [15]	74.81	86.46	73.54	85.34
FaceNet [35]	66.3	82.3	-	-
FNM [32]	80.4	91.2	-	-
PR-REM [4]	83.4	92.1	-	-
PF-cpGAN [38]	86.1	93.8	84.21	90.02
Lin <i>et al.</i> [20]	89.85	<b>95.99</b>	87.55	<b>95.08</b>
SSA [19]	<b>90.91</b>	95.90	88.27	94.88
ours	90.05	95.70	<b>88.35</b>	94.87

The **IJB-B** [49] is challenging in the wild dataset, which was further extended to **IJB-C** [24]. These datasets are of the most challenging benchmarks for FR, containing large pose variation and diversity in resolutions. For evaluating on these datasets, we follow the protocol in [49, 24]. Baseline for comparison on IJB-B and IJB-C datasets are PR-REM [4], FNM [32], FaceNet [35], GOTs [49, 24], VGG-CNN [49, 24], PF-cpGAN [38], CFR-GAN [15], SSA [19], and Lin *et al.* [20]. Table 3 shows that coupled-encoder improves the True Acceptance Rate (TAR) at False Acceptance Rate (FAR) of 0.001. Coupled-encoder performs in

Table 4. Recognition accuracy (%) of the proposed framework on setting1 of Multi-PIE dataset.

loss	view					
	$\pm 15^\circ$	$\pm 30^\circ$	$\pm 45^\circ$	$\pm 60^\circ$	$\pm 75^\circ$	$\pm 90^\circ$
Joint Softmax	99.89	99.71	98.64	94.47	86.92	75.18
MCL	99.82	99.80	98.84	96.17	89.04	77.71
PAC	100	99.87	99.28	98.65	96.49	91.06
PAC+PADA	100	100	99.94	99.71	97.66	93.50
PACM+PADA	100	100	99.97	99.74	98.04	94.64

Table 5. Recognition accuracy (%) of the proposed framework on the Setting1 for Multi-PIE dataset with varying number of negative samples.

Num Negative	view					
	$\pm 15^\circ$	$\pm 30^\circ$	$\pm 45^\circ$	$\pm 60^\circ$	$\pm 75^\circ$	$\pm 90^\circ$
32	100	100	99.91	99.69	97.67	93.19
100	100	100	99.97	99.56	97.97	94.12
150	100	100	99.95	99.94	97.98	94.39
200	100	100	99.97	99.74	98.04	94.64
500	100	100	100	99.61	98.01	94.61

pare with SSA [19] and [20].

### 4.3. Ablation studies

In this section, we analyze the effects of different terms of our proposed framework. First, we study the impact of the loss function’s components, *i.e.*, pose-aware contrastive, memory buffer, and pose-aware adversarial domain adaptation on the performance. Then we report the results of a single encoder. We report on the Multi-PIE dataset to better illustrate the influence of shift in distribution caused by variation in view angle, see Fig. 4. The encoders are pre-trained ResNet50 with Softmax on the VggFace2 dataset. We finetune them with the corresponding loss for the equal number of iterations on the training set of Setting1. The optimizer is SGD for all the experiments. We consider multiple learning rates for training MCL and softmax loss. The learning rate for other experiments is chosen to be 0.001. Also, we have to choose an appropriate margin for training with MCL; therefore, we conduct experiments with different margin values and best results are reported in Table 4.

From Table 4, all the losses perform similarly for the absolute view angles less than  $60^\circ$ . In our experiments, utilizing PAC outperforms Softmax and MCL by almost 15% in identification accuracy for the extreme poses of  $\pm 90^\circ$ . At the same time, it improves the accuracy in near frontal views, except for the  $\pm 45^\circ$ . In the cases of  $\pm 30^\circ$  and  $\pm 45^\circ$ , we observe that when adding the adversarial domain adaptation loss, the performance is improved by 0.13% and 0.66%, respectively. This emphasizes the role of asymmetric mapping in aligning profile and frontal representations.

Contrastive learning frameworks benefit from a larger

number of negative samples [44]. Thus, we expect improvement by integrating the memory bank into the loss function, consistent with our experiment. In Table 5, we further study the effect of varying the number of negative samples on the Setting1 of Multi-PIE dataset. The performance constantly improves until 200 negative samples and then saturates. This saturation is due to the fact that, in the Setting1 of the Multi-PIE, only 150 identities are used for training, and most of the 500 negative samples represent repetitive identities. Moreover, Table 2 shows the performance of the proposed framework with single and coupled-encoder. As the pose disagreement between probe and gallery images increases, coupled encoder presents more improvement, which emphasizes that we need a dedicated frontal encoder to have more flexible mapping for profile faces.

## 5. Conclusion

In this paper, we focused on solving FR in extreme pose scenario. We proposed a new coupled-encoder framework with two distinct encoders that maximize the mutual information between the embeddings of profile and frontal face images. For this goal, we adopt a pose-aware contrastive loss and pose-aware asymmetric training. They force the coupled-encoder to map faces with the same identity to close representations and faces with different identities to the far representations. Furthermore, the memory buffer improves the effectiveness of suggested contrastive learning, by looking at a massive number of identities compared to the mini-batch size. We conducted experiments on multiple benchmarks, showing the capability of our approach to outperform SOTA methods. These performance improvements illustrate the effect of a domain-dedicated feature extractor and employing PACM loss on projecting images to an embedding space where all the images of the same person are close together and far from other individuals, regardless of the view angle. Moreover, the role of each part of our loss function is investigated in the ablation study.

## 6. Acknowledgement

This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. 2022-21102100001. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

## References

- [1] P. Aghdaie, B. Chaudhary, S. Soleymani, J. Dawson, and N. M. Nasrabadi. Morph detection enhanced by structured group sparsity. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 311–320, 2022.
- [2] T. Ahonen, A. Hadid, and M. Pietikainen. Face description with local binary patterns: Application to face recognition. *IEEE Transactions on pattern analysis and machine intelligence*, 28(12):2037–2041, 2006.
- [3] P. Bachman, R. D. Hjelm, and W. Buchwalter. Learning representations by maximizing mutual information across views. *Advances in Neural Information Processing Systems*, 32, 2019.
- [4] K. Cao, Y. Rong, C. Li, X. Tang, and C. C. Loy. Pose-robust face recognition via deep residual equivariant mapping. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5187–5196, 2018.
- [5] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. VggFace2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 67–74. IEEE, 2018.
- [6] J.-C. Chen, V. M. Patel, and R. Chellappa. Unconstrained face verification using deep CNN features. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–9. IEEE, 2016.
- [7] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019.
- [8] H. Du, H. Shi, Y. Liu, J. Wang, Z. Lei, D. Zeng, and T. Mei. Semi-siamese training for shallow face learning. In *European Conference on Computer Vision*, pages 36–53. Springer, 2020.
- [9] G. Elsayed, D. Krishnan, H. Mobahi, K. Regan, and S. Bengio. Large margin deep networks for classification. *Advances in Neural Information Processing Systems*, 31, 2018.
- [10] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. *Image and vision computing*, 28(5):807–813, 2010.
- [11] G. Hinton, O. Vinyals, J. Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- [12] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin. Learning a unified classifier incrementally via rebalancing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 831–839, 2019.
- [13] Y. Hu, X. Wu, B. Yu, R. He, and Z. Sun. Pose-guided photorealistic face rotation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8398–8406, 2018.
- [14] R. Huang, S. Zhang, T. Li, and R. He. Beyond face rotation: Global and local perception GAN for photorealistic and identity preserving frontal view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2439–2448, 2017.
- [15] Y.-J. Ju, G.-H. Lee, J.-H. Hong, and S.-W. Lee. Complete face recovery GAN: Unsupervised joint face rotation and de-occlusion from a single-view image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3711–3721, 2022.
- [16] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- [17] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020.
- [18] X. Li, F. Wang, Q. Hu, and C. Leng. Airface: Lightweight and efficient model for face recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [19] C.-H. Lin and B.-F. Wu. Domain adapting ability of self-supervised learning for face recognition. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 479–483. IEEE, 2021.
- [20] C.-H. Lin and B.-F. Wu. Mitigating domain mismatch in face recognition using style matching. *Neurocomputing*, 487:9–21, 2022.
- [21] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. Sphereface: Deep hypersphere embedding for face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 212–220, 2017.
- [22] 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 Transactions on pattern analysis and machine intelligence*, 41(2):379–393, 2018.
- [23] I. Masi, Y. Wu, T. Hassner, and P. Natarajan. Deep face recognition: A survey. In *2018 31st SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*, pages 471–478. IEEE, 2018.
- [24] B. Maze, J. Adams, J. A. Duncan, N. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson, J. Cheney, et al. Iarpa janus benchmark-c: Face dataset and protocol. In *2018 International Conference on Biometrics (ICB)*, pages 158–165. IEEE, 2018.
- [25] B. Meden, P. Rot, P. Terhörst, N. Damer, A. Kuijper, W. J. Scheirer, A. Ross, P. Peer, and V. Štruc. Privacy-enhancing face biometrics: A comprehensive survey. *IEEE Transactions on Information Forensics and Security*, 2021.
- [26] Q. Meng, X. Xu, X. Wang, Y. Qian, Y. Qin, Z. Wang, C. Zhao, F. Zhou, and Z. Lei. Poseface: Pose-invariant features and pose-adaptive loss for face recognition. *arXiv preprint arXiv:2107.11721*, 2021.
- [27] S. Mosharafian, S. Afzali, Y. Bao, and J. M. Velni. A deep reinforcement learning-based sliding mode control design for partially-known nonlinear systems. *arXiv preprint arXiv:2205.02975*, 2022.
- [28] S. Mosharafian, M. Razzaghpour, Y. P. Fallah, and J. M. Velni. Gaussian process based stochastic model predictive control for cooperative adaptive cruise control. In *2021*

- IEEE Vehicular Networking Conference (VNC)*, pages 17–23. IEEE, 2021.
- [29] M. Nourelahi, F. Dadboud, H. Khalili, A. Niakan, and H. Parsaei. A machine learning model for predicting favorable outcome in severe traumatic brain injury patients after 6 months. *Acute and critical care*, 37(1):45–52, 2022.
- [30] M. Nourelahi, L. Kotthoff, P. Chen, and A. Nguyen. How explainable are adversarially-robust cnns? *arXiv preprint arXiv:2205.13042*, 2022.
- [31] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [32] Y. Qian, W. Deng, and J. Hu. Unsupervised face normalization with extreme pose and expression in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9851–9858, 2019.
- [33] N. Ruiz, E. Chong, and J. M. Rehg. Fine-grained head pose estimation without keypoints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2074–2083, 2018.
- [34] M. Saffari, M. Williams, M. Khodayar, M. Shafie-khah, and J. P. Catalão. Robust wind speed forecasting: A deep spatio-temporal approach. In *2021 IEEE International Conference on Environment and Electrical Engineering and 2021 IEEE Industrial and Commercial Power Systems Europe (EEEIC/I&CPS Europe)*, pages 1–6. IEEE, 2021.
- [35] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015.
- [36] S. Sengupta, J.-C. Chen, C. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs. Frontal to profile face verification in the wild. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–9. IEEE, 2016.
- [37] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. *Advances in Neural Information Processing Systems*, 27, 2014.
- [38] F. Taherkhani, V. Talreja, J. Dawson, M. C. Valenti, and N. M. Nasrabadi. PF-cpGAN: Profile to frontal coupled gan for face recognition in the wild. In *2020 IEEE International Joint Conference on Biometrics (IJCB)*, pages 1–10. IEEE, 2020.
- [39] Y. Tian, D. Krishnan, and P. Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019.
- [40] L. Tran, X. Yin, and X. Liu. Disentangled representation learning GAN for pose-invariant face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1415–1424, 2017.
- [41] D. J. Trosten, S. Lokse, R. Jenssen, and M. Kampffmeyer. Reconsidering representation alignment for multi-view clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1255–1265, 2021.
- [42] X. Tu, J. Zhao, Q. Liu, W. Ai, G. Guo, Z. Li, W. Liu, and J. Feng. Joint face image restoration and frontalization for recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [43] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7167–7176, 2017.
- [44] A. Van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv e-prints*, pages arXiv–1807, 2018.
- [45] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille. Normface: L2 hypersphere embedding for face verification. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1041–1049, 2017.
- [46] Y. Wang, D. Gong, Z. Zhou, X. Ji, H. Wang, Z. Li, W. Liu, and T. Zhang. Orthogonal deep features decomposition for age-invariant face recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 738–753, 2018.
- [47] Y. Wei, M. Liu, H. Wang, R. Zhu, G. Hu, and W. Zuo. Learning flow-based feature warping for face frontalization with illumination inconsistent supervision. In *European Conference on Computer Vision*, pages 558–574. Springer, 2020.
- [48] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *European Conference on Computer Vision*, pages 499–515. Springer, 2016.
- [49] C. Whitelam, E. Taborsky, A. Blanton, B. Maze, J. Adams, T. Miller, N. Kalka, A. K. Jain, J. A. Duncan, K. Allen, et al. Iarpa janus benchmark-b face dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition workshops*, pages 90–98, 2017.
- [50] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.
- [51] X. Yin and X. Liu. Multi-task convolutional neural network for pose-invariant face recognition. *IEEE Transactions on Image Processing*, 27(2):964–975, 2017.
- [52] X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker. Towards large-pose face frontalization in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3990–3999, 2017.
- [53] Y. You, I. Gitman, and B. Ginsburg. Scaling sgd batch size to 32k for imagenet training. *arXiv preprint arXiv:1708.03888*, 6(12):6, 2017.
- [54] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016.
- [55] Z. Zhang and M. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in Neural Information Processing Systems*, 31, 2018.
- [56] J. Zhao, Y. Cheng, Y. Xu, L. Xiong, J. Li, F. Zhao, K. Jayashree, S. Pranata, S. Shen, J. Xing, et al. Towards pose invariant face recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2207–2216, 2018.