

Multi-Prototype Networks for Unconstrained Set-based Face Recognition

Jian Zhao^{1,2,3*}, Jianshu Li¹, Xiaoguang Tu^{1,4}, Fang Zhao⁵, Yuan Xin³, Junliang Xing⁶, Hengzhu Liu², Shuicheng Yan^{1,7}, Jiashi Feng¹

¹National University of Singapore, ²National University of Defense Technology

³Tencent FiT DeepSea AI Lab, ⁴University of Electronic Science and Technology of China

⁵Inception Institute of Artificial Intelligence, ⁶Institute of Automation, Chinese Academy of Sciences, ⁷Qihoo 360 AI Institute

{zhaojian90, jianshu}@u.nus.edu, xguangtu@outlook.com, zhaofang0627@gmail.com, maxxin@tencent.com

jlxing@nlpr.ia.ac.cn, hengzhuliu@nudt.edu.cn, {eleyans, elefjia}@nus.edu.sg

Abstract

In this paper, we study the challenging unconstrained set-based face recognition problem where each subject face is instantiated by a set of media (images and videos) instead of a single image. Naively aggregating information from all the media within a set would suffer from the large intra-set variance caused by heterogeneous factors (e.g., varying media modalities, poses and illuminations) and fail to learn discriminative face representations. A novel **Multi-Prototype Network (MPNet)** model is thus proposed to learn multiple prototype face representations adaptively from the media sets. Each learned prototype is representative for the subject face under certain condition in terms of pose, illumination and media modality. Instead of handcrafting the set partition for prototype learning, MPNet introduces a **Dense SubGraph (DSG)** learning sub-net that implicitly untangles inconsistent media and learns a number of representative prototypes. Qualitative and quantitative experiments clearly demonstrate superiority of the proposed model over state-of-the-arts.

1. Introduction

Recent advances of deep learning approaches have remarkably boosted the performance of face recognition. Some approaches claim to have achieved [36, 4, 19, 48, 46] or even surpassed [31, 40, 47] human performance on several benchmarks. However, those approaches only recognize faces over a single image or video sequence. Such scenarios deviate from the reality. In practical face recognition systems (and arguably human cortex for face recognition), each subject face to recognize is often enrolled with

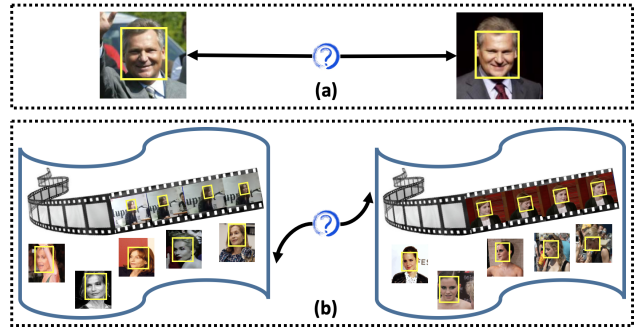


Figure 1: Difference between (a) face recognition over a single image per subject and (b) unconstrained set-based face recognition. For unconstrained set-based face recognition, each subject is represented by a set of mixed images and videos captured under unconstrained conditions. Each set contains large variations in face pose, expression, illumination and occlusion issues. Existing single-medium based recognition approaches cannot successfully address this problem. Best viewed in color.

a set of images and videos captured under varying conditions and acquisition methods. Intuitively such rich information can benefit face recognition performance, which however has not been effectively exploited in existing approaches [45, 21, 48].

In this paper, we consider the challenging task—*unconstrained set-based face recognition* firstly introduced in [18]—that is more consistent with real-world scenarios. Unconstrained set-based face recognition defines the minimal facial representation unit as a set of images and videos instead of a single medium. Set-based face recognition requires solving a more difficult set-to-set matching problem, where both the probe and gallery are sets of face media. This task raises the necessity to build *subject-specific* face models for each subject individually, instead of relying on a single multi-class recognition model as before. An illustration on the difference between traditional face recognition over a single input image and

*Jian Zhao is the corresponding author. Homepage: <https://zhaoj9014.github.io>. This work was done during Jian Zhao served as a short-term “Texpert” research scientist at Tencent FiT DeepSea AI Lab, Shenzhen, China.

the targeted face recognition over a set of unconstrained images/videos is given in Fig. 1. The most significant challenge in the unconstrained set-based face recognition task is how to learn good representations for the media set, even in presence of *large intra-set variance* of real-world subject faces caused by varying conditions in illumination, sensor, compression, *etc.*, and subject attributes such as facial pose, expression and occlusion. Solving this problem needs to address these distracting factors effectively and learn set-level discriminative face representation.

Recently, several set-based face recognition methods have been proposed [3, 5, 7, 13]. They generally adopt the following two strategies to obtain set-level face representation. One is to learn a set of image-level face representations from each face medium in the set individually [7, 23], and use all the information for following face recognition. Such a strategy is obviously computationally expensive as it needs to perform exhaustive pairwise matching and is fragile to outlier medium captured under unusual conditions. The other strategy is to aggregate face representations within the set through simple average or max pooling and generate single representation for each set [5, 30]. However, this obviously suffers from information loss and insufficient exploitation of the image/video set.

To overcome the limitations of existing methods, we propose a novel **Multi-Prototype Network** (MPNet) model. To learn better set-level representations, MPNet introduces a **Dense SubGraph** (DSG) learning sub-net to implicitly factorize each face media set of a particular subject into a number of disentangled sub-sets, instead of handcrafting the set partition using some intuitive features. Each dense sub-graph discovers a sub-set (representing a *prototype*) of face media that are with small intra-set variance but discriminative from other subject faces. MPNet learns to enhance the compactness of the prototypes as well as their coverage of large variance for a single subject face, through which heterogeneous attributes within each face media set are sufficiently considered and flexibly untangled. This significantly helps improve the unconstrained set-based face recognition performance by providing multiple comprehensive and succinct face representations, reducing the impact of media inconsistency. Compared with existing set-based face recognition methods [3, 5, 7, 13], MPNet effectively addresses the large variance challenge and offers more discriminative and flexible face representations with lower computational complexity. Also, superior to naive average or max pooling of face features, MPNet effectively preserves the necessary information through the DSG learning for set-based face recognition. The main contributions of this work can be summarized as follows:

- We propose a novel and effective multi-prototype discriminative learning architecture MPNet. To our best knowledge, MPNet is the first end-to-end trainable

model that adaptively learns multiple prototype face representations from sets of media. It is effective at addressing the large intra-set variance issue that is critical to set-based face recognition.

- MPNet introduces a **Dense SubGraph** (DSG) learning sub-net that automatically factorizes each face media set into a number of disentangled prototypes representing consistent face media with sufficient discriminativeness. Through the DSG sub-net, MPNet is capable of untangling inconsistent media and dealing with faces captured under challenging conditions robustly.
- DSG provides a general loss that encourages compactness around multiple discovered centers with strong discrimination. It offers a new and systematic approach for large variance object recognition in the real world.

Based on the above technical contributions, we have presented a high-performance model for unconstrained set-based face recognition. It achieves currently best results on IJB-A [18], YTF [42] and IJB-C [25] benchmark datasets with significant improvement over state-of-the-arts.

2. Related Work

Recent top performing approaches for face recognition often rely on deep CNNs with advanced architectures. For instance, the VGGface model [27, 2], as an application of the VGG architecture [32], provides state-of-the-art performance. The DeepFace model [36, 37] also uses a deep CNN coupled with 3D alignment. FaceNet [31] utilizes the inception deep CNN architecture for unconstrained face recognition. DeepID2+ [35] and DeepID3 [34] extend the FaceNet model by including joint Bayesian metric learning and multi-task learning, yielding better unconstrained face recognition performance. SphereFace [20], CosFace [40], AM-Softmax [39] and ArcFace [11] exploit margin-based representation learning to achieve small intra-class distance and large inter-class distance. Those methods enhance their overall performance via carefully designed architectures, which are however not tailored for unconstrained set-based face recognition.

With the introduction of IJB-A benchmark [18] by NIST in 2015, the problem of unconstrained set-based face recognition attracts increasing attention. Recent solutions to this problem are also based on deep architectures, which are leading approaches on LFW [16] and YTF [42]. Among them, B-CNN [7] applies a new **Bilinear CNN** (B-CNN) architecture for face identification. Pooling Faces [13] aligns faces in 3D and partitions them according to facial and imaging properties. PAMs [23] handles pose variability by learning **Pose-Aware Models** (PAMs) for frontal, half-profile and full-profile poses to perform face recognition in the wild. Those methods often employ separate processing

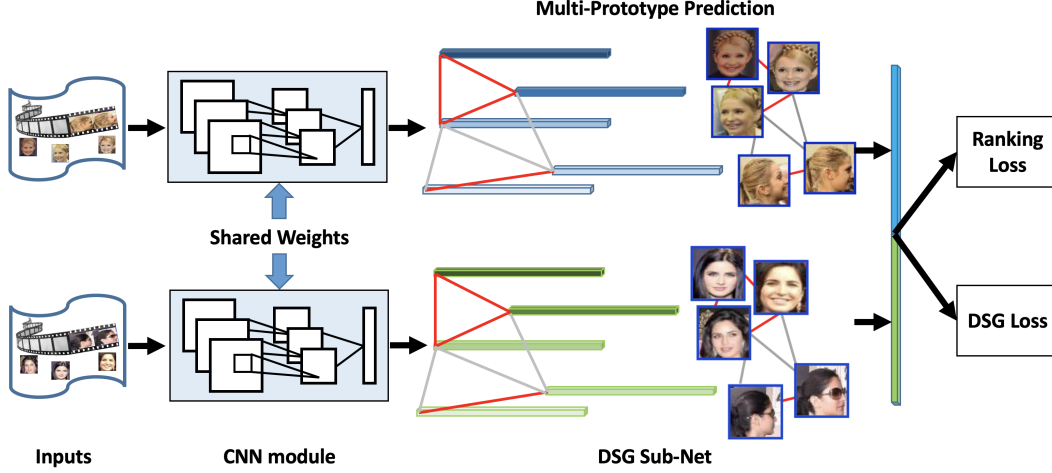


Figure 2: The proposed **Multi-Prototype Network (MPNet)** for unconstrained set-based face recognition. MPNet takes a pair of face media sets as input and outputs a matching result, *i.e.*, same person or not. It adopts a modern deep siamese CNN architecture for deep set-based facial representation learning, and introduces a new DSG sub-net to learn discriminative prototypes for each set. MPNet is end-to-end trainable through the ranking loss and auxiliary DSG loss. Best viewed in color.

steps without considering the modality variance within one set of face media and underlying multiple prototype structures. Therefore, much useful information may loss, leading to inferior performance.

Our proposed MPNet shares a similar idea as subcategory-aware object classification [12] that considers intra-class variance in building object classifiers, and [9, 10] that considers to construct a graph-matching metric for domain adaption. Our method differs from them in following aspects: 1) the “prototype” is not pre-defined in MPNet; 2) learning DSG within each face media set implicitly discovers consistent faces sharing similar conditions, through which heterogeneous factors are flexibly untangled; 3) MPNet is based on deep learning and can be end-to-end trainable. It is also interesting to investigate the application of our MPNet architecture in generic object recognition tasks.

3. Multi-Prototype Networks

Fig. 2 visualizes the architecture of the MPNet, which takes a pair of face media sets as input and outputs a matching result for the unconstrained set-based face recognition. It adopts a modern deep siamese CNN architecture for set-based facial representation learning, and introduces a new DSG sub-net for learning the multi-prototype that models various representative faces under different conditions for the input. MPNet is end-to-end trainable by minimizing the ranking loss and a new DSG loss. We now present each component in detail.

3.1. Set-based Facial Representation Learning

Different from face recognition over a single image, the task of set-based face recognition aims to accept or reject

the claimed identity of a subject represented by a face media set containing both images and videos. Performance is assessed using two measures: percentage of false accepts and that of false rejects. A good model should optimize both metrics simultaneously. MPNet is designed to nonlinearly map the raw sets of faces to multiple prototypes in a low dimensional space such that the distance between these prototypes is small if the sets belong to the same subject, and large otherwise. The similarity metric learning is achieved by training MPNet with two identical CNN branches that share weights. MPNet handles inputs in a pair-wise, set-to-set way so that it explicitly organizes the face media in a way favorable to set-based face recognition.

MPNet learns face representations at multi-scale for gaining strengthened robustness to scale variance in real-world faces. Specifically, for each medium within a face media set, a multi-scale pyramid is constructed by resizing the image or video frame to four different scales. To handle the error of face detection, MPNet performs random cropping to collect local and global patches from each scale of the multi-scale pyramid with a fixed size, as illustrated in Fig. 3. To handle the imbalance of realistic face data (*e.g.*, some subjects are enrolled with scarce media from limited images while some with redundant media from reduplicative video frames), the data distribution within each set is adjusted by resampling. In particular, the set containing scarce media (*i.e.*, less than a pre-defined parameter R that is set empirically) is augmented by duplicating and flipping images, which is intuitively beneficial with the support from more relevant information. The large set with redundant media (*i.e.*, more than R) is subsampled to the size of R . The resulting input streams to MPNet are tuples of face

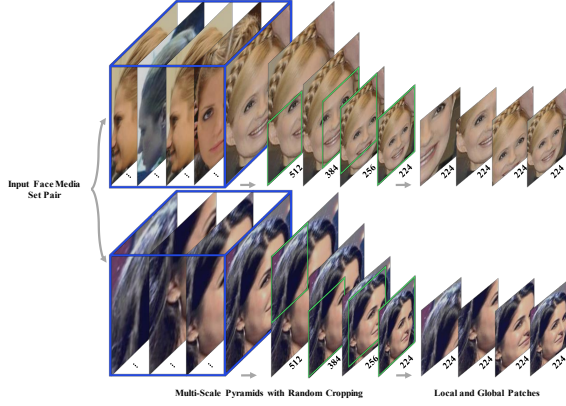


Figure 3: Illustration of multi-scale pyramid construction and random cropping strategy for the MPNet. For each medium within a face media set, a multi-scale pyramid is constructed by resizing the medium to four different scales. Local and global patches are then randomly cropped from each multi-scale pyramid with a fixed size. Best viewed in color.

media set pairs and the associated ground truth annotations $\{(X^{p1}, X^{p2}, y^p)\}$, where X^{p1} and X^{p2} denote the two sets of the p -th pair, and y^p denotes the binary pair-wise label.

The proposed MPNet adopts a siamese CNN architecture in which two branches share weights for pairwise set-based facial representation learning. Each branch is initiated with VGGface [27], including 13 convolutional layers, 5 pooling layers and 2 fully-connected layers. We make the following careful architectural design for each branch to ensure that the learned deep facial representations are more suitable for multi-prototype representation learning. 1) For activation functions, instead of using ReLU [26] to suppress all the negative responses, we adopt the PReLU [14] to allow negative responses. PReLU improves model fitting with nearly zero extra computational cost and little overfitting risk, benefiting convergence of MPNet. 2) We adopt two local normalization layers after the 2^{nd} convolutional layer and the 4^{th} convolutional layer, respectively. The local normalization tends to uniformize the mean and variance of a feature map around a local neighborhood. This is especially useful for correcting non-uniform illumination or shading artifacts. 3) We adopt an average operator for the last pooling layer and a max operator for the previous 4 pooling layers to generate compact and discriminative deep set-based facial representations. Note that our approach is not restricted to the CNN module used, and can also be generalized to other state-of-the-art architectures for performance boosting.

The learned deep facial representation for each face media set is denoted as $\{f_1, f_2, \dots, f_R\}$. Here recall R is the specified size of the face media set after distribution balance.

3.2. Multi-Prototype Discriminative Learning

Throughout this paper, a prototype is defined as a collection of similar face media that are representative for a subject face under certain conditions. Face media forming the same prototype have small variance and one can safely extract representation by pooling approaches without worrying about information loss.

To address the critical large variance issue in set-based face recognition, we propose the multi-prototype discriminative learning. With this component, each face media set is implicitly factorized into a certain number of prototypes. Multi-prototype learning encourages the output facial representations to be compact w.r.t. certain prototypes and meanwhile discriminative for different subjects. Thus, MPNet is capable of modelling the prototype-level interactions effectively while addressing the large variance and false matching caused by untypical faces. MPNet dynamically offers an optimal trade-off between facial information preserving and computation cost. It does not require exhaustive matching across all of the possible pairs from two sets for accurately recognizing faces in the wild. It learns multiple prototypes through a dense subgraph learning as detailed below.

To discover the underlying multiple prototypes of each face media set instead of handcrafting the set partition, we propose a novel DSG learning approach. DSG formulates the similarity of face media within a set through a graph and discovers the prototype by mining the dense subgraphs. Each subgraph has high internal similarity and small similarity to the outside media. Compared with clustering-based data partition, DSG is advantageous in terms of flexibility and robustness to outliers. Each subgraph provides a prototype for the input subject faces. We then perform face recognition at the prototype level, which is concise and also sufficiently informative.

Given a latent affinity graph characterizing the relations among entities (face media here), denoted as $\mathcal{G}=(\mathcal{V}, \mathcal{E})$, each dense subgraph refers to a prototype of the vertices (\mathcal{V}) with larger internal connection (\mathcal{E}) affinity than other candidates. In this work, learning DSG within each face media set implicitly discovers consistent faces sharing similar conditions such as age, pose, expression, illumination and media modality, through which heterogeneous factors are flexibly untangled.

Formally, suppose the graph \mathcal{G} is associated with an affinity A , and each element of A encodes the similarity between two face media: $a_{ij} = \text{aff}(f_i, f_j)$. Let K be the number of prototypes (or equivalently, number of dense subgraphs) and $Z = [z_1, \dots, z_K] \in \mathbb{R}^{n \times K}$ be the partition indicator: $z_{ik} = 1$ indicates the i -th medium is in the k -th prototype. The DSG aims to find the representative

subgraph via optimizing Z through

$$\max_Z \text{tr}(Z^\top AZ), \text{ s.t. }, z_{ik} \in \{0, 1\}, Z\mathbf{1} = \mathbf{1}, \quad (1)$$

where $\mathbf{1}$ is an all-1 vector. Here the 2^{nd} constraint guarantees that every media will be allocated to only one prototype. The allocation after learning would maximize the intra-prototype media similarity. This is significantly different from k-means clustering where the centers are not necessary to learn and similarity is not defined based on the distance to the center.

This problem is not easy to solve. We therefore carefully design the DSG layers that form the DSG sub-net to optimize it end-to-end. This sub-net consists of two layers, which takes the set-based facial representations f_i 's as input and outputs the reconstructed discriminative features.

The 1^{st} layer makes the prototype prediction. Given the input deep facial representation f_i , this layer outputs its indicator $z_i \in \{0, 1\}$ by dynamically projecting each input latent affinity graph to K prototypes, $z_i = \sigma(W^\top f_i)$, where σ is the sigmoid activation function to rectify inputs to $[0, 1]$ and W is the DSG predictor parameter. Given the predicted z_i and input f_i , the 2^{nd} layer computes the DSG loss (see Sec. 3.3) to ensure the reconstructed representation \hat{f}_i form reasonably compact and discriminative prototypes. \hat{f}_i is obtained by element-wisely multiplying the output of the 2^{nd} layer of the DSG sub-net with the predicted multi-prototype indicator Z from its 1^{st} layer. More details are given in Sec. 4.

3.3. Optimization

We optimize MPNet by minimizing the following two loss functions in a conjugate way.

Ranking loss: a ranking loss is designed in MPNet to enforce the distance to shrink for genuine set pairs and be large for imposter set pairs, so that MPNet explicitly maps input patterns into the target spaces to approximate the semantic distance in the input space.

We first ℓ_2 -normalize the outputs from the DSG sub-net, so that all the set-based facial representations are within the same range for loss computation. Then, we use Euclidean distance to measure the fine-grained pairwise dissimilarity between \hat{f}_i and \hat{f}_j :

$$d_{ij} = \|\hat{f}_i - \hat{f}_j\|_2^2, \quad (2)$$

where $i, j \in \{1, \dots, R\}$.

We further ensemble the R^2 distances into one energy-based matching result for each coarse-level set pair:

$$E = \frac{\sum_{i,j}^R d_{ij} \times \exp(\beta d_{ij})}{\sum_{i,j}^R \exp(\beta d_{ij})}, \quad (3)$$

where β is a bandwidth parameter.

Our final ranking loss function is formulated as

$$\mathcal{L}_{\text{Ranking}}(\hat{f}_i) \triangleq \min_E \{(1 - y^p)E + y^p \max(0, \tau - E)\}, \quad (4)$$

where τ is a margin, such that $E_G + \tau < E_I$, E_G is the distance for genuine pair, E_I is the distance for imposter pair, and y^p is the binary pairwise label, with 0 for genuine pair ($\mathcal{L}_{\text{Ranking}} = E_G$ in Eq. (4)) and 1 for imposter pair ($\mathcal{L}_{\text{Ranking}} = \max(0, \tau - E_I)$ in Eq. (4)).

Dense SubGraph loss: We propose to learn dense prototypes through solving the problem defined in Eq. (1). Expanding the objective in Eq. (1) gives

$$\text{tr}(Z^\top AZ) = \sum_{i,j=1}^n \sum_{k=1}^K z_{ik} a_{ij} z_{jk}. \quad (5)$$

Since $z_{ik}, z_{jk} \in \{0, 1\}$, we have $\sum_{k=1}^K z_{ik} a_{ij} z_{jk} = a_{ij}$ only if $z_{ik} = z_{jk} = 1$ for some k , i.e., the face media i and j are divided into the same prototype. Maximizing the trace in Eq. (1) is to find the partition of samples (indicated by z) to form subgraphs such that the samples associated with the same subgraph have the largest total affinity (i.e., density) $\sum_{z_{ik} \neq 0, z_{jk} \neq 0} a_{ij}$. In practice, maximizing $\text{tr}(Z^\top AZ)$ would encourage contributions of the representations \hat{f}_i belonging to the same prototype to be close to each other and each resulted cluster to be far away from others, i.e., they form multiple dense subgraphs.

In Eq. (5), each element of the affinity A encodes similarity between two corresponding media, i.e.

$$a_{ij} \triangleq \exp\left\{\frac{-d_{ij}}{\delta^2}\right\}. \quad (6)$$

Equivalently, the DSG learning can be achieved through the following minimization problem:

$$\min_Z \text{tr}(Z^\top DZ), \text{ s.t. } z_{ij} \in \{0, 1\}, Z\mathbf{1} = \mathbf{1}, \quad (7)$$

where D is the Euclidean distance matrix computed in Eqn. (2).

Then we define the following DSG loss function to optimize the learned deep set-based facial representations:

$$\mathcal{L}_{\text{DSG}}(\hat{f}_i) \triangleq \left\{ \min_Z \text{tr}(Z^\top DZ), \text{ s.t. } z_{ij} \in \{0, 1\}, Z\mathbf{1} = \mathbf{1} \right\}. \quad (8)$$

Thus, minimizing the DSG loss would encourage contributions of the representations \hat{f}_i 's belonging to the same prototype to be close to each other. If one visualizes the learned representations in the high-dimensional space, the learned representations of one face media set form several *compact* clusters and each cluster may be far away from others. In this way, a face media set with large variance is distributed to several clusters implicitly. Each cluster has a

Method	Verification			Identification		
	TAR@FAR=0.10	TAR@FAR=0.01	TAR@FAR=0.001	FNIR@FPIR=0.10	FNIR@FPIR=0.01	Rank1
OpenBR [18]	0.433±0.006	0.236±0.009	0.104±0.014	0.851±0.028	0.934±0.017	0.246±0.011
GOTS [18]	0.627±0.012	0.406±0.014	0.198±0.008	0.765±0.033	0.953±0.024	0.433±0.021
BCNNs [7]	-	-	-	0.659±0.032	0.857±0.024	0.588±0.020
Pooling faces [13]	0.631	0.309	-	-	-	0.846
LSFS [38]	0.895±0.013	0.733±0.034	0.514±0.060	0.387±0.032	0.617±0.063	0.820±0.024
Deep Multi-pose [1]	0.991	0.787	-	0.250	0.48	0.846
DCNN _{manual} +metric [6]	0.947±0.011	0.787±0.043	-	-	-	0.852±0.018
Triplet Similarity [30]	0.945±0.002	0.790±0.030	0.590±0.050	0.246±0.014	0.444±0.065	0.880±0.015
PAMs [23]	-	0.826±0.018	0.652±0.037	-	-	0.840±0.012
DCNN _{fusion} [5]	0.967±0.009	0.838±0.042	-	0.210±0.033	0.423±0.094	0.903±0.012
Masi <i>et al.</i> [24]	-	0.886	0.725	-	-	0.906
Triplet Embedding [30]	0.964±0.005	0.900±0.010	0.813±0.002	0.137±0.014	0.247±0.030	0.932±0.010
All-In-One [29]	0.976±0.004	0.922±0.010	0.823±0.020	0.113±0.014	0.208±0.020	0.947±0.008
Template Adaptation [8]	0.979±0.004	0.939±0.013	0.836±0.027	0.118±0.016	0.226±0.049	0.928±0.001
NAN [44]	0.979±0.004	0.941±0.008	0.881±0.011	0.083±0.009	0.183±0.041	0.958±0.005
DA-GAN [48]	0.991±0.003	0.976±0.007	0.930±0.005	0.051±0.009	0.110±0.039	0.971±0.007
ℓ_2 -softmax [28]	0.984±0.002	0.970±0.004	0.943±0.005	0.044±0.006	0.085±0.041	0.973±0.005
3D-PIM [47]	0.996±0.001	0.989±0.002	0.977±0.004	0.016±0.005	0.064±0.045	0.990±0.002
baseline	0.968±0.009	0.871±0.014	0.735±0.031	0.188±0.011	0.372±0.045	0.907±0.010
w/o DSG	0.971±0.006	0.887±0.012	0.743±0.027	0.182±0.010	0.367±0.041	0.912±0.008
MPNet _{K=3}	0.971±0.006	0.863±0.019	0.734±0.033	0.189±0.013	0.386±0.043	0.909±0.007
MPNet _{K=10}	0.971±0.007	0.880±0.015	0.740±0.026	0.179±0.009	0.361±0.044	0.913±0.009
MPNet _{K=200}	0.979±0.004	0.924±0.013	0.764±0.022	0.171±0.012	0.350±0.046	0.923±0.008
MPNet _{K=500}	0.980±0.005	0.919±0.013	0.779±0.021	0.169±0.009	0.337±0.042	0.932±0.008
MPNet _{K=1000}	0.975±0.008	0.909±0.017	0.757±0.025	0.164±0.011	0.359±0.040	0.926±0.010
MPNet (Ours)	0.997±0.002	0.991±0.003	0.984±0.005	0.011±0.005	0.059±0.040	0.994±0.003

Table 1: Face recognition performance comparison on IJB-A. The results are averaged over 10 testing splits. “-” means the result is not reported. Standard deviation is not available for some methods.

Method	Acc
DeepID [33]	0.932
DeepFace [36]	0.914
Center loss [41]	0.949
SphereFace [20]	0.950
FaceNet [31]	0.951
VGGface [27]	0.973
CosFace [40]	0.976
ArcFace [11]	0.980
MPNet (Ours)	0.991

Table 2: Face recognition performance comparison on YTF.

small variance. We also conduct experiments for illustration in Sec. 4.1.

To simplify the above optimization, we propose to relax the constraint of $z_{ij} \in \{0, 1\}$ to $0 \leq z_{ij} \leq 1$ by a sigmoid activation function. Thus, the DSG loss is re-defined as

$$\mathcal{L}_{\text{DSG}}(\hat{f}_i) \triangleq \left\{ \min_Z \text{tr}(Z^\top DZ), \text{ s.t. } 0 \leq z_{ij} \leq 1, Z\mathbf{1} = \mathbf{1} \right\}. \quad (9)$$

We adopt the joint supervision of ranking and DSG losses to train MPNet for multi-prototype discriminative learning:

$$\mathcal{L} = \mathcal{L}_{\text{Ranking}} + \lambda \mathcal{L}_{\text{DSG}}, \quad (10)$$

where λ is a weighting parameter among the two losses.

Clearly, MPNet is end-to-end trainable and can be optimized with BP and SGD algorithm. We summarize the learning algorithm of MPNet in Algorithm 1.

4. Experiments

We evaluate MPNet qualitatively and quantitatively under various settings for unconstrained set-based face recognition on IJB-A [18], YTF [42] and IJB-C [25].

Implementation Details We initialize the CNN module of MPNet for deep set-based facial representation learning

with the VGGface model [27], and fine-tune it on the target dataset. For each medium with the provided face bounding box, we first crop the facial RoI accordingly and then resize it to multiple $r \times r \times 3$ sizes to build the multi-scale pyramids, where $r=224, 256, 384$ and 512 . The size of inputs to MPNet is fixed as $224 \times 224 \times 3$ by randomly cropping local and global patches of compatible size from images/video frames. No 2D or 3D face alignment is used. The threshold R for balancing input data distribution is set as 128 for trading-off recognition accuracy and computation cost. The weights of the 1st layer (implemented with a 1D convolution layer with sigmoid activation function) of the DSG subnet are initialized by normal distribution with an std 0.001. The number of total prototypes K is set as 500. We also conduct experiments to illustrate how the K influences the overall performance in Sec. 4.2. The bandwidth parameter β in Eq. (3) is set to 10, the margin τ of the ranking loss is fixed as 0.8, and the trade-off parameter λ is set as 0.01 by 5-fold cross-validation. Different values of λ lead to different deep feature distributions. With proper λ , the discriminative power of deep features can be significantly enhanced. $\lambda = 0.01$ is large enough for balancing the scales of two loss terms as the sub-graph loss calculates summations over more pairs. The proposed network is implemented based on the publicly available Caffe platform [17], which is trained on three NVIDIA GeForce GTX TITAN X GPUs with 12G memory. During training, the learning rate is initialized to 0.01, and during fine-tuning, the learning rate is initialized to 0.001. We train our model using SGD with a batch size of 1 face media set pair, momentum of 0.9, and weight decay of 0.0005.

4.1. Evaluations on IJB-A Benchmark

IJB-A contains 5,397 images and 2,042 videos from 500 subjects, captured from in-the-wild environment to avoid

Algorithm 1 Multi-prototype learning algorithm

Input: Training data $X_p = \{(X^{p1}, X^{p2}, y^p)\}$. Initialized parameters θ, W in the CNN module and DSG sub-net, respectively. Hyperparameters $R, K, \beta, \tau, \lambda$ and learning rate μ^t . The number of iteration $t \leftarrow 0$;

Output: The parameters θ and W .

while not converge **do**

$t \leftarrow t + 1$;

 Compute the joint loss by $\mathcal{L}^t = \mathcal{L}_{Ranking}^t + \lambda \mathcal{L}_{DSG}^t$;

 Compute the backpropagation error for each p by $\frac{\partial \mathcal{L}^t}{\partial X_p^t} = \frac{\partial \mathcal{L}_{Ranking}^t}{\partial X_p^t} + \lambda \cdot \frac{\partial \mathcal{L}_{DSG}^t}{\partial X_p^t}$;

 Update the parameters θ by $\theta^{t+1} = \theta^t - \mu^t \sum_p^m \frac{\partial \mathcal{L}^t}{\partial X_p^t} \cdot \frac{\partial X_p^t}{\partial \theta^t}$;

 Update the parameters W by $W^{t+1} = W^t - \mu^t \sum_p^m \frac{\partial \mathcal{L}^t}{\partial X_p^t} \cdot \frac{\partial X_p^t}{\partial W^t}$;

end while

near frontal bias. For training and testing, 10 random splits are provided by each protocol, respectively. It contains two tasks, face verification and identification. The performance is evaluated by TAR@FAR, FNIR@FPIR and Rank metrics, respectively.

4.1.1 Ablation Study and Quantitative Comparison

We first investigate different architectures and losses of MPNet to see their respective roles in unconstrained set-based face recognition. We compare 8 variants of MPNet, *i.e.*, baseline (siamese VGGface [27]), w/o DSG, MPNet_{K $\in\{3,10,200,500,1000\}$} , and MPNet (backbone: ResNet-101 [15]).

The performance comparison in terms of TAR@FAR, FNIR@FPIR and Rank1 on IJB-A is reported in the lower panel of Tab. 1. By comparing the results from w/o DSG vs. the baseline, around 1% improvement for overall evaluation metrics can be observed. This confirms the benefits of the basic refining tricks in terms of the network structure. Compared with w/o DSG, MPNet_{K=500} further boosts the performance by around 3%, which speaks well for the superiority of using the auxiliary DSG loss to enhance the deep set-based facial representation learning. It simplifies unconstrained set-based face recognition, yet reserves discriminative and comprehensive information. By varying the numbers of prototypes, one can see that as K increases from 3 to 1,000, the performance on the overall metrics improves consistently when $K \leq 500$. This demonstrates that the affinity-based dense subgraph learning of the proposed DSG sub-net can effectively enhance the deep feature capacity of unconstrained set-based face recognition. However, further increasing K does not bring further performance improvement and may even harm the performance on the overall metrics. The reason is that an appropriately large value of K will predict a sparse prototype partition indicator matrix Z , which helps reach an optimal trade-off between facial information preserving and computation cost

for addressing the large variance and false matching caused by untypical faces. However, an oversize value of K will enforce the learned filters to all zero ones, which always produces invariant performance without any discriminative information. We hence set K to 500 in all the experiments.

For fair comparison with other state-of-the-arts (upper panel of Tab. 1), we further replace the backbone from VG-Gface to ResNet-101 (bottom row) while keeping other settings the same. Our MPNet achieves the best results over 10 testing splits on both protocols. This superior performance demonstrates that MPNet is very effective for the unconstrained set-based face recognition in presence of large intra-set variance. Compared with existing set-based face recognition approaches, our MPNet can effectively address the large variance challenge and offer more discriminative and flexible face representations with small computational complexity. Also, superior to the naive average or max pooling of face features, MPNet effectively preserves necessary information through the DSG learning for set-based face recognition.

Moreover, compared with exhaustive matching strategies (*e.g.*, DCNN [5]) which have $O(mn)$ complexity for similarity computation (m, n are media numbers of each face set to recognize) and take $\sim 1.2s$ for recognizing each probe set, our MPNet is more efficient as it operates on prototype level, which significantly reduces the computational complexity to $O(K^2)$, $K^2 \ll mn$ (K is the prototype number) and takes $\sim 0.5s$ for recognizing each probe set. Although naive average or max pooling strategies (*e.g.*, Pooling faces [13]) are slightly advantageous in testing time ($\sim 0.3s$ for recognizing each probe set), they suffer from information loss severely. Our MPNet effectively preserves the necessary information through DSG learning for unconstrained set-based face recognition.

4.1.2 Qualitative Comparison

We then verify the effectiveness of our deep multi-prototype discriminative learning strategy. The predicted prototypes with relatively larger affinities within the set 1311 and set 3038 from the testing data of IJB-A split1 are visualized using t-SNE [22] in Fig. 4. We observe that MPNet explicitly learns to automatically predict the prototype memberships within each coarse-level face set reflecting different poses (*e.g.*, the first 6 learned prototypes), expressions (*e.g.*, the 1st and the 7th learned prototypes), illumination (*e.g.*, the 2nd and 5th learned prototypes), and media modalities (*e.g.*, the 5th and 6th learned prototypes). Each learned prototype contains coherent media offering collective facial representation with specific patterns. Outliers within each face set are detected by MPNet (*e.g.*, the last learned prototypes). MPNet is learnt to enhance the compactness of the prototypes as well as their coverage of large variance for a single subject face, through which the heterogeneous attributes within each face media set are sufficiently considered and flexibly untangled. Compared with clustering-based data partition, MPNet with DSG learning is advantageous since it is end-to-end trainable, can learn more discriminative features and is robust to outliers. Learning DSG maximizes the intra-prototype media similarity and inter-prototype difference, resulting in discriminative face representations. This is significantly different from clustering (*e.g.*, k-means) methods where only the similarity defined based on the distance to the center is considered during learning.

Finally, we visualize the verification results in Fig. 5 for IJB-A split1 to gain insight into unconstrained set-based face recognition. After computing the similarities for all pairs of probe and reference sets, we sort the resulting list. Each row represents a probe and reference set pair. The original sets within IJB-A contain from one to dozens of media. Up to 8 individual media are shown with the last space showing a mosaic of the remaining media in the set. Between the sets are the set IDs for probe and reference as well as the best matched and best non-matched similarities. Fig. 5 (blue, left) shows the best matched cases. In the top-30 scoring correct matches, we immediately note that every reference set contains dozens of media. The probe sets either contain dozens of media or one medium that matches well. Fig. 5 (blue, right) shows the worst matched cases, representing failed matching. The thirty lowest matched results from single-medium probe sets are all under extremely challenging unconstrained conditions. These extremely difficult cases cannot be solved even using the specific operations designed in MPNet. Fig. 5 (green, left) showing the worst non-matched cases highlights the understandable errors involving single-medium probe sets representing impostors in challenging orientations. Fig. 5 (green, right) showing the best non-matched cases shows the most cer-

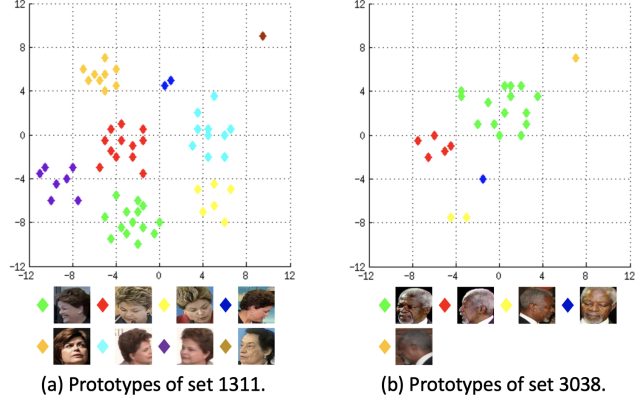


Figure 4: Visualization of learned prototypes within set 1311 (a) and set 3038 (b) by MPNet, from the testing data of IJB-A split1. Each colored cluster shows a learned prototype. One sampled face within each prototype is shown for better illustration. Best viewed in color.

tain non-mates, again often involving large sets with enough guidance from the relevant information of the same subject.

4.2. Evaluations on YTF Benchmark

YTF contains 3,425 videos of 1,595 different subjects. The average length of a video clip is 181.3 frames. All the video sequences were downloaded from YouTube. We follow the unrestricted with labeled outside data protocol and report the result on 5,000 video pairs.

The face recognition performance comparison of the proposed MPNet with other state-of-the-arts on YTF is reported in Tab 2. MPNet improves the 2nd-best by 1.10%, which well verified the superiority of MPNet for effectively learning set-level discriminative face representations.

4.3. Evaluations on IJB-C Benchmark

IJB-C contains 31,334 images and 11,779 videos from 3,531 subjects, which are split into 117,542 frames, 8.87 images and 3.34 videos per subject, captured from in-the-wild environments to avoid the near frontal bias. For fair comparison, we follow the template-based setting and evaluate models on the standard 1:1 verification protocol in terms of TAR@FAR.

The face recognition performance comparison of the proposed MPNet with other state-of-the-arts on IJB-C is reported in Tab. 3. MPNet beats the 2nd-best by 5.60% w.r.t. TAR@FAR=10⁻⁵, which further shows its remarkable generalizability for recognizing faces in the wild, and the learned deep features are robust and disambiguated.

5. Conclusion

We proposed a novel **Multi-Prototype Network (MPNet)** with a new **Dense SubGraph (DSG)** learning sub-net to ad-



Figure 5: Verification results analysis for IJB-A split1. (blue, left) The best matched cases, (blue, right) The worst matched cases, (green, left) The worst non-matched cases (green, right) The best non-matched cases. For better viewing of this figure, please see the original zoomed-in color pdf file.

Method	$\text{TAR@FAR}=10^{-5}$	$\text{TAR@FAR}=10^{-4}$	$\text{TAR@FAR}=10^{-3}$	$\text{TAR@FAR}=10^{-2}$
GOTS [25]	0.066	0.147	0.330	0.620
FaceNet [31]	0.330	0.487	0.665	0.817
VGGface [27]	0.437	0.598	0.748	0.871
VGGface2.ft [2]	0.768	0.862	0.927	0.967
MN-vc [43]	0.771	0.862	0.927	0.968
MPNet	0.827	0.898	0.940	0.971

Table 3: Face recognition performance comparison on IJB-C.

dress unconstrained set-based face recognition, which adaptively learns compact and discriminative multi-prototype representations. Comprehensive experiments demonstrate the superiority of MPNet over state-of-the-arts. The proposed framework can be easily extended to other generic object recognition tasks by utilizing the area-specific sets. In future, we will explore a pure MPNet architecture where all components are replaced with well designed MPNet layers, which can hierarchically exploit the multi-prototype discriminative information to solve complex computer vision problems.

Acknowledgement

The work of Jian Zhao was partially supported by China Scholarship Council (CSC) grant 201503170248.

The work of Junliang Xing was partially supported by the National Science Foundation of China 61672519.

The work of Jiashi Feng was partially supported by NUS IDS R-263-000-C67-646, ECRA R-263-000-C87-133 and MOE Tier-II R-263-000-D17-112.

References

- [1] W. AbdAlmageed, Y. Wu, S. Rawls, S. Harel, T. Hassner, I. Masi, J. Choi, J. Lekust, J. Kim, P. Natarajan, et al. Face recognition using deep multi-pose representations. In *WACV*, pages 1–9, 2016.
- [2] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. VggFace2: a dataset for recognising faces across pose and age. In *FG*, pages 67–74, 2018.
- [3] R. Chellappa, J.-C. Chen, R. Ranjan, S. Sankaranarayanan, A. Kumar, V. M. Patel, and C. D. Castillo. Towards the design of an end-to-end automated system for image and video-based recognition. *arXiv preprint arXiv:1601.07883*, 2016.
- [4] J. Chen, V. M. Patel, L. Liu, V. Kellokumpu, G. Zhao, M. Pietikäinen, and R. Chellappa. Robust local features for remote face recognition. *IVC*, 64:34–46, 2017.
- [5] J.-C. Chen, V. M. Patel, and R. Chellappa. Unconstrained face verification using deep cnn features. In *WACV*, pages 1–9, 2016.
- [6] J.-C. Chen, R. Ranjan, A. Kumar, C.-H. Chen, V. M. Patel, and R. Chellappa. An end-to-end system for unconstrained face verification with deep convolutional neural networks. In *ICCVW*, pages 118–126, 2015.
- [7] A. R. Chowdhury, T.-Y. Lin, S. Maji, and E. Learned-Miller. One-to-many face recognition with bilinear CNNs. In *WACV*, pages 1–9, 2016.
- [8] N. Crosswhite, J. Byrne, C. Stauffer, O. Parkhi, Q. Cao, and A. Zisserman. Template adaptation for face verification and identification. In *FG*, pages 1–8, 2017.
- [9] D. Das and C. G. Lee. Sample-to-sample correspondence for unsupervised domain adaptation. *EAAI*, 73:80–91, 2018.
- [10] D. Das and C. G. Lee. Unsupervised domain adaptation using regularized hyper-graph matching. In *ICIP*, pages 3758–3762, 2018.

- [11] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. ArcFace: additive angular margin loss for deep face recognition. *arXiv preprint arXiv:1801.07698*, 2018.
- [12] J. Dong, W. Xia, Q. Chen, J. Feng, Z. Huang, and S. Yan. Subcategory-aware object classification. In *CVPR*, pages 827–834, 2013.
- [13] T. Hassner, I. Masi, J. Kim, J. Choi, S. Harel, P. Natarajan, and G. Medioni. Pooling faces: template based face recognition with pooled face images. In *CVPRW*, pages 59–67, 2016.
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: surpassing human-level performance on imagenet classification. In *ICCV*, pages 1026–1034, 2015.
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [16] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: a database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: convolutional architecture for fast feature embedding. In *ACM MM*, pages 675–678, 2014.
- [18] B. F. Klare, B. Klein, E. Taborsky, A. Blanton, J. Cheney, K. Allen, P. Grother, A. Mah, M. Burge, and A. K. Jain. Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus Benchmark A. In *CVPR*, pages 1931–1939, 2015.
- [19] J. Li, J. Zhao, F. Zhao, H. Liu, J. Li, S. Shen, J. Feng, and T. Sim. Robust face recognition with deep multi-view representation learning. In *ACM MM*, pages 1068–1072, 2016.
- [20] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. SphereFace: Deep hypersphere embedding for face recognition. In *CVPR*, volume 1, page 1, 2017.
- [21] X. Lu, Y. Wang, W. Zhang, S. Ding, and W. Jiang. Deep CNNs for face verification. In *CCBR*, pages 85–92, 2016.
- [22] L. v. d. Maaten and G. Hinton. Visualizing data using t-SNE. *JMLR*, 9(Nov):2579–2605, 2008.
- [23] I. Masi, S. Rawls, G. Medioni, and P. Natarajan. Pose-aware face recognition in the wild. In *CVPR*, pages 4838–4846, 2016.
- [24] I. Masi, A. T. Tran, J. T. Leksut, T. Hassner, and G. Medioni. Do we really need to collect millions of faces for effective face recognition? *arXiv preprint arXiv:1603.07057*, 2016.
- [25] 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 *ICB*, 2018.
- [26] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In *ICML*, pages 807–814, 2010.
- [27] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015.
- [28] R. Ranjan, C. D. Castillo, and R. Chellappa. L2-constrained Softmax loss for discriminative face verification. *arXiv preprint arXiv:1703.09507*, 2017.
- [29] R. Ranjan, S. Sankaranarayanan, C. D. Castillo, and R. Chellappa. An all-in-one convolutional neural network for face analysis. *arXiv preprint arXiv:1611.00851*, 2016.
- [30] S. Sankaranarayanan, A. Alavi, C. D. Castillo, and R. Chellappa. Triplet probabilistic embedding for face verification and clustering. In *BTAS*, pages 1–8, 2016.
- [31] F. Schroff, D. Kalenichenko, and J. Philbin. FaceNet: a unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015.
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [33] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *NIPS*, pages 1988–1996, 2014.
- [34] Y. Sun, D. Liang, X. Wang, and X. Tang. DeepID3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015.
- [35] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. In *CVPR*, pages 2892–2900, 2015.
- [36] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. DeepFace: closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014.
- [37] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Web-scale training for face identification. In *CVPR*, pages 2746–2754, 2015.
- [38] D. Wang, C. Otto, and A. K. Jain. Face search at scale: 80 million gallery. *arXiv preprint arXiv:1507.07242*, 2015.
- [39] F. Wang, J. Cheng, W. Liu, and H. Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018.
- [40] H. Wang, Y. Wang, Z. Zhou, X. Ji, and W. Liu. CosFace: large margin cosine loss for deep face recognition. In *CVPR*, pages 5265–5274, 2018.
- [41] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *ECCV*, pages 499–515, 2016.
- [42] L. Wolf, T. Hassner, and I. Maoz. Face recognition in unconstrained videos with matched background similarity. In *CVPR*, pages 529–534, 2011.
- [43] W. Xie and A. Zisserman. Multicolumn networks for face recognition. *arXiv preprint arXiv:1807.09192*, 2018.
- [44] J. Yang, P. Ren, D. Zhang, D. Chen, F. Wen, H. Li, and G. Hua. Neural aggregation network for video face recognition. In *CVPR*, pages 5216–5225, 2017.
- [45] H. Ye, W. Shao, H. Wang, J. Ma, L. Wang, Y. Zheng, and X. Xue. Face recognition via active annotation and learning. In *ACM MM*, pages 1058–1062, 2016.
- [46] 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 *CVPR*, pages 2207–2216, 2018.
- [47] J. Zhao, L. Xiong, Y. Cheng, Y. Cheng, J. Li, L. Zhou, Y. Xu, J. Karlekar, S. Pranata, S. Shen, J. Xing, S. Yan, and J. Feng. 3d-aided deep pose-invariant face recognition. In *IJCAI*, pages 1184–1190, 2018.

- [48] J. Zhao, L. Xiong, P. K. Jayashree, J. Li, F. Zhao, Z. Wang, P. S. Pranata, P. S. Shen, S. Yan, and J. Feng. Dual-agent GANs for photorealistic and identity preserving profile face synthesis. In *NIPS*, pages 66–76, 2017.