

# Unconstrained Face Verification using Deep CNN Features

Jun-Cheng Chen<sup>1</sup>, Vishal M. Patel<sup>2</sup>, and Rama Chellappa<sup>1</sup>

1. University of Maryland, College Park

2. Rutgers, The State University of New Jersey

pullpull@cs.umd.edu, vishal.m.patel@rutgers.edu, rama@umiacs.umd.edu

## Abstract

*In this paper, we present an algorithm for unconstrained face verification based on deep convolutional features and evaluate it on the newly released IARPA Janus Benchmark A (IJB-A) dataset as well as on the traditional Labeled Face in the Wild (LFW) dataset. The IJB-A dataset includes real-world unconstrained faces from 500 subjects with full pose and illumination variations which are much harder than the LFW and Youtube Face (YTF) datasets. The deep convolutional neural network (DCNN) is trained using the CASIA-WebFace dataset. Results of experimental evaluations on the IJB-A and the LFW datasets are provided.*

## 1. Introduction

Face verification is one of the core problems in computer vision and has been actively researched for over two decades [40]. In face verification, given two videos or images, the objective is to determine whether they belong to the same person. Many algorithms have been shown to work well on images that are collected in controlled settings. However, the performance of these algorithms often degrades significantly on images that have large variations in pose, illumination, expression, aging, cosmetics, and occlusion.

To deal with this problem, many methods have focused on learning invariant and discriminative representation from face images and videos. One approach is to extract over-complete and high-dimensional feature representation followed by a learned metric to project the feature vector into a low-dimensional space and to compute the similarity score. For instance, the high-dimensional multi-scale Local Binary Pattern (LBP)[5] features extracted from local patches around facial landmarks is reasonably effective for face recognition. Face representation based on Fisher vector (FV) has also shown to be effective for face recognition problems [26][23], [9]. However, deep convolutional neural networks (DCNN) have demonstrated impressive performances on different tasks such as object recog-

nition [21][31], object detection [14], and face verification [25]. It has been shown that a DCNN model can not only characterize large data variations but also learn a compact and discriminative feature representation when the size of the training data is sufficiently large. Once the model is learned, it is possible to generalize it to other tasks by fine-tuning the learned model on target datasets [13]. In this work, we train a DCNN model using a relatively small face dataset, the CASIA-WebFace [38], and compare the performance of our method with other commercial off-the-shelf face matchers on the challenging IJB-A dataset which contains significant variations in pose, illumination, expression, resolution and occlusion. We also evaluate the performance of the proposed method on the LFW dataset.

The rest of the paper is organized as follows. We briefly review some related works in Section 2. Details of the different components of the proposed method including the DCNN representation and joint Bayesian metric learning are given in Section 3. The protocol and the experimental results are presented in Section 4. Finally, we conclude the paper in Section 5 with a brief summary and discussion.

## 2. Related Work

In this section, we briefly review several recent related works on face verification.

### 2.1. Feature Learning

Learning invariant and discriminative feature representation is the first step for a face verification system. It can be broadly divided into two categories: (1) hand-crafted features, and (2) feature representation learned from data. In the first category, Ahonen *et al.* [1] showed that the Local Binary Pattern (LBP) is effective for face recognition. Gabor wavelets [39][37] have also been widely used to encode multi-scale and multi-orientation information for face images. Chen *et al.* [6] demonstrated good results for face verification using the high-dimensional multi-scale LBP features extracted from patches around facial landmarks. In the second category, Patel *et al.* [24] and Chen *et al.* [11][10] applied dictionary-based approaches for im-

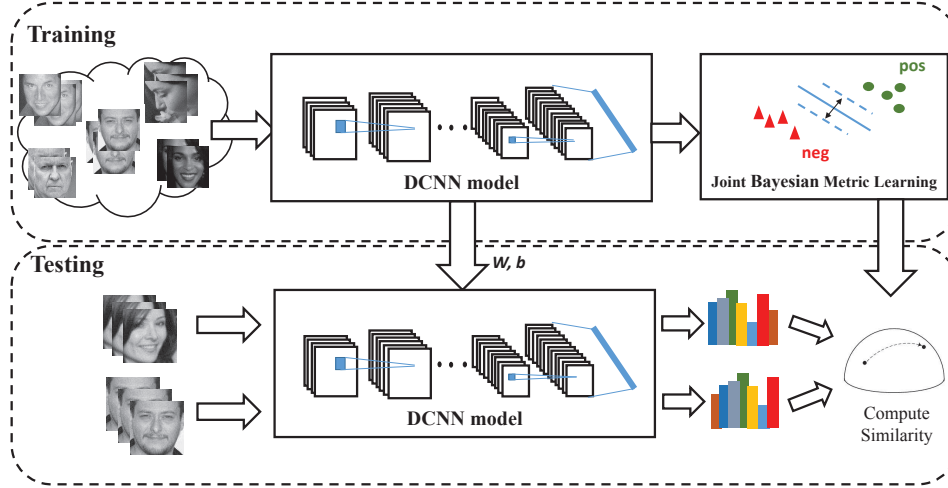


Figure 1. An overview of the proposed DCNN approach for face verification.

age and video-based face recognition by learning representative atoms from the data which are compact and robust to pose and illumination variations. [26][23][7] used the FV encoding to generate over-complete and high-dimensional feature representation for still and video-based face recognition. Lu *et al.* [22] proposed a dictionary learning framework in which the sparse codes of local patches generated from local patch dictionaries are pooled to generate a high-dimensional feature vector. The high-dimensionality of feature vectors makes these methods hard to train and scale to large datasets. However, advances in deep learning methods have shown that compact and discriminative representation can be learned using DCNN from very large datasets. Taigman *et al.* [33] learned a DCNN model on the frontalized faces generated with a general 3D shape model from a large-scale face dataset and achieved better performance than many traditional face verification methods. Sun *et al.* [28][30] achieved results that surpass human performance for face verification on the LFW dataset using an ensemble of 25 simple DCNN with fewer layers trained on weakly aligned face images from a much smaller dataset than the former. Schroff *et al.* [25] adapted the state-of-the-art deep architecture for object recognition to face recognition and trained it on a large-scale unaligned private face dataset with the triplet loss. This method also achieved top performances on face verification problems. These works essentially demonstrate the effectiveness of the DCNN model for feature learning and detection/recognition/verification problems.

## 2.2. Metric Learning

Learning a similarity measure from data is the other key component that can boost the performance of a face verification system. Many approaches have been proposed in

the literature that essentially exploit the label information from face images or face pairs. For instance, Weinberger *et al.* [36] proposed Large Margin Nearest Neighbor (LMNN) metric which enforces the large margin constraint among all triplets of labeled training data. Taigman *et al.* [32] learned the Mahalanobis distance using the Information Theoretic Metric Learning (ITML) method [12]. Chen *et al.* [5] proposed a joint Bayesian approach for face verification which models the joint distribution of a pair of face images instead of the difference between them, and the ratio of between-class and within-class probabilities is used as the similarity measure. Hu *et al.* [17] learned a discriminative metric within the deep neural network framework. Huang *et al.* [18] learned a projection metric over a set of labeled images which preserves the underlying manifold structure.

## 3. Method

Our approach consists of both training and testing stages. For training, we first perform face and landmark detection on the CASIA-WebFace, and the IJB-A datasets to localize and align each face. Next, we train our DCNN on the CASIA-WebFace and derive the joint Bayesian metric using the training sets of the IJB-A dataset and the DCNN features. Then, given a pair of test image sets, we compute the similarity score based on their DCNN features and the learned metric. Figure 1 gives an overview of our method. The details of each component of our approach are presented in the following subsections.

### 3.1. Preprocessing

Before training the convolutional network, we perform landmark detection using the method presented in [2][3] because of its ability to be effective on unconstrained faces. Then, each face is aligned into the canonical coordinate

with similarity transform using the 7 landmark points (*i.e.* two left eye corners, two right eye corners, nose tip, and two mouth corners). After alignment, the face image resolution is  $100 \times 100$  pixels, and the distance between the centers of two eyes is about 36 pixels.

### 3.2. Deep Face Feature Representation

A DCNN with small filters and very deep architecture (*i.e.* 19 layers in [27] and 22 layers in [31]) has shown to produce state-of-the-art results on many datasets including ImageNet 2014, LFW, and Youtube Face dataset. Stacking small filters to approximate large filters and to build very deep convolution networks not only reduces the number of parameters but also increases the nonlinearity of the network. In addition, the resulting feature representation is compact and discriminative.

Our approach is motivated by [38]. However, we only consider the identity information per face without modeling the pair-wise cost. The dimensionality of the input layer is  $100 \times 100 \times 1$  for gray-scale images. The network includes 10 convolutional layers, 5 pooling layers and 1 fully connected layer. The detailed architecture is shown in Table 1. Each convolutional layer is followed by a rectified linear unit (ReLU) except the last one, Conv52. Instead of suppressing all the negative responses to zero using ReLU, we use parametric ReLU (PReLU)[16] which allows negative responses that in turn improves the network performance. Thus, we use PReLU as an alternative to ReLU in our work. Moreover, two local normalization layers are added after Conv12 and Conv22, respectively to mitigate the effect of illumination variations. The kernel size of all filters is  $3 \times 3$ . The first four pooling layers use the max operator. To generate a compact and discriminative feature representation, we use average pooling for the last layer, pool<sub>5</sub>. The feature dimensionality of pool<sub>5</sub> is thus equal to the number of channel of Conv52 which is 320. Dropout ratio is set as 0.4 to regularize Fc6 due to the large number of parameters (*i.e.*  $320 \times 10548$ ). To classify a large number of subjects in the training data (*i.e.* 10548), this low-dimensional feature should contain strong discriminative information from all the face images. Consequently, the pool<sub>5</sub> feature is used for face representation. The extracted features are further  $L_2$ -normalized into unit length before the metric learning stage. If there are multiple frames available for the subject, we use the average of the pool<sub>5</sub> features as the overall feature representation. Figure 2 illustrates some of the extracted feature maps.

### 3.3. Joint Bayesian Metric Learning

To utilize the positive and negative label information available from the training dataset, we learn a joint Bayesian metric which has achieved good performances on face verification problems [5][4]. Instead of modeling the differ-

ence vector between two faces, this approach directly models the joint distribution of feature vectors of both  $i$ th and  $j$ th images,  $\{\mathbf{x}_i, \mathbf{x}_j\}$ , as a Gaussian. Let  $P(\mathbf{x}_i, \mathbf{x}_j|H_I) \sim N(0, \Sigma_I)$  when  $\mathbf{x}_i$  and  $\mathbf{x}_j$  belong to the same class, and  $P(\mathbf{x}_i, \mathbf{x}_j|H_E) \sim N(0, \Sigma_E)$  when they are from different classes. In addition, each face vector can be modeled as,  $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\mu}$  stands for the identity and  $\boldsymbol{\epsilon}$  for pose, illumination, and other variations. Both  $\boldsymbol{\mu}$  and  $\boldsymbol{\epsilon}$  are assumed to be independent zero-mean Gaussian distributions,  $N(0, \mathbf{S}_\mu)$  and  $N(0, \mathbf{S}_\epsilon)$ , respectively.

The log likelihood ratio of intra- and inter-classes,  $r(\mathbf{x}_i, \mathbf{x}_j)$ , can be computed as follows:

$$r(\mathbf{x}_i, \mathbf{x}_j) = \log \frac{P(\mathbf{x}_i, \mathbf{x}_j|H_I)}{P(\mathbf{x}_i, \mathbf{x}_j|H_E)} = \mathbf{x}_i^T \mathbf{M} \mathbf{x}_i + \mathbf{x}_j^T \mathbf{M} \mathbf{x}_j - 2\mathbf{x}_i^T \mathbf{R} \mathbf{x}_j, \quad (1)$$

where  $\mathbf{M}$  and  $\mathbf{R}$  are both negative semi-definite matrices. Equation (1) can be rewritten as  $(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j) - 2\mathbf{x}_i^T \mathbf{B} \mathbf{x}_j$  where  $\mathbf{B} = \mathbf{R} - \mathbf{M}$ . More details can be found in [5]. Instead of using the EM algorithm to estimate  $\mathbf{S}_\mu$  and  $\mathbf{S}_\epsilon$ , we optimize the distance in a large-margin framework as follows:

$$\underset{\mathbf{M}, \mathbf{B}, b}{\operatorname{argmin}} \sum_{i,j} \max[1 - y_{ij}(b - (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j) + 2\mathbf{x}_i^T \mathbf{B} \mathbf{x}_j), 0], \quad (2)$$

where  $b \in \mathbb{R}$  is the threshold, and  $y_{ij}$  is the label of a pair:  $y_{ij} = 1$  if person  $i$  and  $j$  are the same and  $y_{ij} = -1$ , otherwise. For simplicity, we denote  $(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j) - 2\mathbf{x}_i^T \mathbf{B} \mathbf{x}_j$  as  $d_{\mathbf{M}, \mathbf{B}}(\mathbf{x}_i, \mathbf{x}_j)$ .  $\mathbf{M}$  and  $\mathbf{B}$  are updated using stochastic gradient descent as follows and are equally trained on positive and negative pairs in turn:

$$\begin{aligned} \mathbf{M}_{t+1} &= \begin{cases} \mathbf{M}_t, & \text{if } y_{ij}(b_t - d_{\mathbf{M}, \mathbf{B}}(\mathbf{x}_i, \mathbf{x}_j)) > 1 \\ \mathbf{M}_t - \gamma y_{ij} \boldsymbol{\Gamma}_{ij}, & \text{otherwise,} \end{cases} \\ \mathbf{B}_{t+1} &= \begin{cases} \mathbf{B}_t, & \text{if } y_{ij}(b_t - d_{\mathbf{M}, \mathbf{B}}(\mathbf{x}_i, \mathbf{x}_j)) > 1 \\ \mathbf{B}_t + 2\gamma y_{ij} \mathbf{x}_i \mathbf{x}_j^T, & \text{otherwise,} \end{cases} \\ b_{t+1} &= \begin{cases} b_t, & \text{if } y_{ij}(b_t - d_{\mathbf{M}, \mathbf{B}}(\mathbf{x}_i, \mathbf{x}_j)) > 1 \\ b_t + \gamma_b y_{ij}, & \text{otherwise,} \end{cases} \end{aligned} \quad (3)$$

where  $\boldsymbol{\Gamma}_{ij} = (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T$  and  $\gamma$  is the learning rate for  $\mathbf{M}$  and  $\mathbf{B}$ , and  $\gamma_b$  for the bias  $b$ . We use random semi-definite matrices to initialize both  $\mathbf{M} = \mathbf{V}\mathbf{V}^T$  and  $\mathbf{B} = \mathbf{W}\mathbf{W}^T$  where both  $\mathbf{V}$  and  $\mathbf{W} \in \mathbb{R}^{d \times d}$ , and  $v_{ij}$  and  $w_{ij} \sim N(0, 1)$ . Note that  $\mathbf{M}$  and  $\mathbf{B}$  are updated only when the constraints are violated. In our implementation, the ratio of the positive and negative pairs that we generate based on the identity information of the training set is 1:20. In addition, the other reason to train the metric instead of using traditional EM is that for IJB-A training and test data, some templates only contain a single image. More details about the IJB-A dataset are given in Section 4.

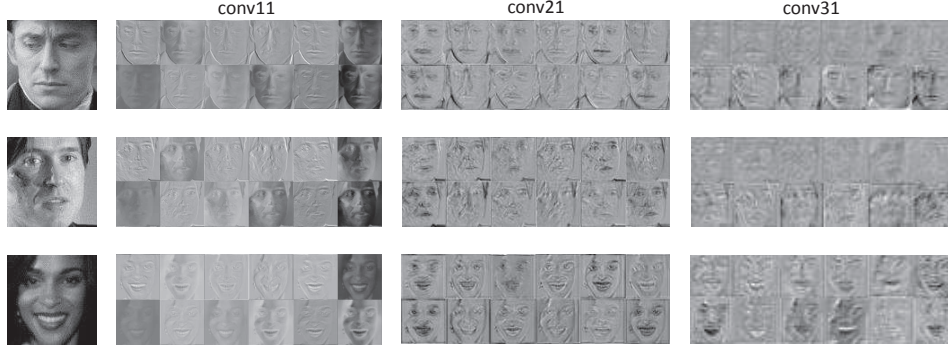


Figure 2. An illustration of some feature maps of Conv11, Conv21, and Conv31 layers. At the upper layers, the feature maps capture more global shape features which are also more robust to illumination changes than Conv11.

| Name    | Type             | Filter Size/Stride          | Output Size                | Depth | #Params |
|---------|------------------|-----------------------------|----------------------------|-------|---------|
| Conv11  | convolution      | $3 \times 3 \times 1 / 1$   | $100 \times 100 \times 32$ | 1     | 0.28K   |
| Conv12  | convolution      | $3 \times 3 \times 32 / 1$  | $100 \times 100 \times 64$ | 1     | 18K     |
| Pool1   | max pooling      | $2 \times 2 / 2$            | $50 \times 50 \times 64$   | 0     |         |
| Conv21  | convolution      | $3 \times 3 \times 64 / 1$  | $50 \times 50 \times 64$   | 1     | 36K     |
| Conv22  | convolution      | $3 \times 3 \times 64 / 1$  | $50 \times 50 \times 128$  | 1     | 72K     |
| Pool2   | max pooling      | $2 \times 2 / 2$            | $25 \times 25 \times 128$  | 0     |         |
| Conv31  | convolution      | $3 \times 3 \times 128 / 1$ | $25 \times 25 \times 96$   | 1     | 108K    |
| Conv32  | convolution      | $3 \times 3 \times 96 / 1$  | $25 \times 25 \times 192$  | 1     | 162K    |
| Pool3   | max pooling      | $2 \times 2 / 2$            | $13 \times 13 \times 192$  | 0     |         |
| Conv41  | convolution      | $3 \times 3 \times 192 / 1$ | $13 \times 13 \times 128$  | 1     | 216K    |
| Conv42  | convolution      | $3 \times 3 \times 128 / 1$ | $13 \times 13 \times 256$  | 1     | 288K    |
| Pool4   | max pooling      | $2 \times 2 / 2$            | $7 \times 7 \times 256$    | 0     |         |
| Conv51  | convolution      | $3 \times 3 \times 256 / 1$ | $7 \times 7 \times 160$    | 1     | 360K    |
| Conv52  | convolution      | $3 \times 3 \times 160 / 1$ | $7 \times 7 \times 320$    | 1     | 450K    |
| Pool5   | avg pooling      | $7 \times 7 / 1$            | $1 \times 1 \times 320$    | 0     |         |
| Dropout | dropout (40%)    |                             | $1 \times 1 \times 320$    | 0     |         |
| Fc6     | fully connection |                             | 10548                      | 1     | 3296K   |
| Cost    | softmax          |                             | 10548                      | 0     |         |
| total   |                  |                             |                            | 11    | 5006K   |

Table 1. The architecture of DCNN used in this paper.

### 3.4. DCNN Training Details

The DCNN is implemented using caffe[19] and trained on the CASIA-WebFace dataset. The CASIA-WebFace dataset contains 494,414 face images of 10,575 subjects downloaded from the IMDB website. After removing the 27 overlapping subjects with the IJB-A dataset, there are 10548 subjects<sup>1</sup> and 490,356 face images. For each subject, there still exists several false images with wrong identity labels and few duplicate images. All images are scaled into  $[0, 1]$  and subtracted from the mean. The data is augmented with horizontal flipped face images. We use the standard batch size 128 for the training phase. Because it only contains sparse positive and negative pairs per batch in addition to the false image problems, we do not take the verification

<sup>1</sup>The list of overlapping subjects is available at [http://www.umi.acs.umd.edu/~pullpull/janus\\_overlap.xlsx](http://www.umi.acs.umd.edu/~pullpull/janus_overlap.xlsx)

cost into consideration as is done in [30]. The initial negative slope for PReLU is set to 0.25 as suggested in [16]. The weight decay of all convolutional layers are set to 0, and the weight decay of the final fully connected layer to  $5e-4$ . In addition, the learning rate is set to  $1e-2$  initially and reduced by half every 100,000 iterations. The momentum is set to 0.9. Finally, we use the snapshot of 1,000,000th iteration for all our experiments.

## 4. Experiments

In this section, we present the results of the proposed approach on the challenging IARPA Janus Benchmark A (IJB-A) [20], its extended version Janus Challenging set 2 (JANUS CS2) dataset and the LFW dataset. The JANUS CS2 dataset contains not only the sampled frames and images in the IJB-A but also the original videos. The JANUS



CS2 dataset<sup>2</sup> includes much more test data for identification and verification problems in the defined protocols than the IJB-A dataset. The receiver operating characteristic curves (ROC) and the cumulative match characteristic (CMC) scores are used to evaluate the performance of different algorithms. The ROC curve measures the performance in the verification scenarios, and the CMC score measures the accuracy in a closed set identification scenarios.

#### 4.1. JANUS-CS2 and IJB-A

Both the IJB-A and JANUS CS2 contain 500 subjects with 5,397 images and 2,042 videos split into 20,412 frames, 11.4 images and 4.2 videos per subject. Sample images and video frames from the datasets are shown in Fig. 3. The videos are only released for the JANUS CS2 dataset. The IJB-A evaluation protocol consists of verification (1:1 matching) over 10 splits. Each split contains around 11,748 pairs of templates (1,756 positive and 9,992 negative pairs) on average. Similarly, the identification (1:N search) protocol also consists of 10 splits which evaluates the search performance. In each search split, there are about 112 gallery templates and 1763 probe templates (*i.e.* 1,187 genuine probe templates and 576 impostor probe templates). On the other hand, for the JANUS CS2, there are about 167 gallery templates and 1763 probe templates and all of them are used for both identification and verification. The training set for both dataset contains 333 subjects, and the test set contains 167 subjects. Ten random splits of training and testing are provided by each benchmark, respectively. The main differences between IJB-A and JANUS CS2 evaluation protocol are (1) IJB-A considers the open-set identification problem and the JANUS CS2 considers the closed-set identification and (2) IJB-A considers the more difficult pairs which are the subsets from the JANUS CS2 dataset.



Figure 3. Sample images and frames from the IJB-A and JANUS CS2 datasets. A variety of challenging variations on pose, illumination, resolution, occlusion, and image quality are present in these images.

Both the IJB-A and the JANUS CS2 datasets are divided into training and test sets. For the test sets of both benchmarks, the image and video frames of each subject are randomly split into gallery and probe sets without any overlapping subjects between them. Unlike the LFW and YTF

<sup>2</sup>The JANUS CS2 dataset is not publicly available yet.

datasets which only use a sparse set of negative pairs to evaluate the verification performance, the IJB-A and JANUS CS2 both divide the images/video frames into gallery and probe sets so that it uses all the available positive and negative pairs for the evaluation. Also, each gallery and probe set consist of multiple templates. Each template contains a combination of images or frames sampled from multiple image sets or videos of a subject. For example, the size of the similarity matrix for JANUS CS2 split1 is  $167 \times 1806$  where 167 are for the gallery set and 1806 for the probe set (*i.e.* the same subject reappears multiple times in different probe templates). Moreover, some templates contain only one profile face with challenging pose with low quality image. In contrast to the LFW and YTF datasets which only include faces detected by the Viola Jones face detector [34], the images in the IJB-A and JANUS CS2 contain extreme pose, illumination and expression variations. These factors essentially make the IJB-A and JANUS CS2 challenging face recognition datasets [20].

#### 4.2. Evaluation on JANUS-CS2 and IJB-A

For the JANUS CS2 dataset, we compare the results of our DCNN method with the FV approach proposed in [26] and two other commercial off-the-shelf matchers, COTS1 and GOTS [20]. The COTS1 and GOTS baselines provided by JANUS CS2 are the top performers from the most recent NIST FRVT study [15]. The FV method is trained on the LFW dataset which contains few faces with extreme pose. Therefore, we use the pose information estimated from the landmark detector and select face images/video frames whose yaw angle are less than or equal to  $\pm 25$  degrees for each gallery and probe set. If there are no images/frames satisfying the constraint, we choose the one closest to the frontal one. However, for the DCNN method, we use all the frames without applying the same selection strategy.<sup>3</sup> Figures 4 and 5 show the ROC curves and the CMC curves, respectively for the verification results using the previously described protocol where DCNN means using DCNN feature with cosine distance, “ft” means finetuning on the training data, “metric” means applying Joint Bayesian metric learning, and “color” means to use all of the RGB images instead of gray-scale images. For the results of  $\text{DCNN}_{ft+metric}$ , besides finetuning and metric learning, we also replace ReLU with PReLU and apply data augmentation (*i.e.* randomly cropping  $100 \times 100$ -pixel subregions from a  $125 \times 125$  region). For  $\text{DCNN}_{ft+metric+color}$ <sup>4</sup>, we further use RGB images and larger face regions. (*i.e.* we use  $125 \times 125$ -pixel face regions and resize them into  $100 \times 100$ -pixel ones.) Then, we show the fusion results,

<sup>3</sup>We fix the typos in [8] that the selection strategy is only applied to FV-based method, not for DCNN.

<sup>4</sup> $\text{DCNN}_{ft+metric+color}$  and  $\text{DCNN}_{fusion}$  are our improved results for JANUS CS2 and IJB-A datasets obtained after the paper was accepted.



















| Probe Template  | Rank-1   | Rank-2   | Rank-3   | Rank-4  | Rank-5  |
|---|--|--|--|---|---|
| #Image: 22<br><br>Template ID: 2047<br>Subject ID: 543 | #Image: 14<br><br><b>Template ID: 2030</b><br><b>Subject ID: 543</b>  | #Image: 3<br><br>Template ID: 5794<br>Subject ID:: 791  | #Image: 34<br><br>Template ID: 226<br>Subject ID: 102   | #Image: 32<br><br>Template ID: 187<br>Subject ID: 101  | #Image: 50<br><br>Template ID: 4726<br>Subject ID: 404 |
| #Image: 1<br><br>Template ID: 2993<br>Subject ID: 1559 | #Image: 22<br><br><b>Template ID: 2992</b><br><b>Subject ID: 1559</b> | #Image: 9<br><br>Template ID: 948<br>Subject ID:: 1558  | #Image: 6<br><br>Template ID: 1312<br>Subject ID:: 1704 | #Image: 4<br><br>Template ID: 3779<br>Subject ID: 1876 | #Image: 4<br><br>Template ID: 5812<br>Subject ID: 2166 |
| #Image: 1<br><br>Template ID: 2062<br>Subject ID: 158  | #Image: 25<br><br>Template ID: 986<br>Subject ID: 347                 | #Image: 7<br><br>Template ID: 5295<br>Subject ID:: 2058 | #Image: 3<br><br>Template ID: 3729<br>Subject ID: 606   | #Image: 32<br><br>Template ID: 187<br>Subject ID: 101  | #Image: 6<br><br>Template ID: 5494<br>Subject ID: 2102 |

Table 2. Query results. The first column shows the query images from probe templates. The remaining 5 columns show the corresponding top-5 queried gallery templates.

DCNN<sub>fusion</sub>, by directly summing the similarity scores of two models, DCNN<sub>ft+metric</sub> and DCNN<sub>ft+metric+color</sub>, where DCNN<sub>ft+metric</sub> is trained on gray-scale images with smaller face regions and DCNN<sub>ft+metric+color</sub> is trained on RGB images with larger face regions. From these figures, we can clearly see the impact of each component to the improvement of final identification and verification results. From the ROC and CMC curves, we see that the DCNN method performs better than other competitive methods. This can be attributed to the fact that the DCNN model does capture face variations over a large dataset and generalizes well to a new small dataset.

We illustrate the query samples in Table 2. The first column shows the query images from the probe templates. The remaining five columns show the corresponding top-5 queried gallery templates (*i.e.* rank-1 means the most similar one, rank-2 the second most similar, etc.). For the first two rows, our approach can successfully find the subjects in rank 1. For the third, the query template only contains one image with extreme pose. However, in the corresponding gallery template for the same subject, it happens to contain only near-frontal faces. Thus, it failed to find the subject within the top-5 matches. To solve the pose generalization problem of CNN features, one possible solution is to augment the templates by synthesizing faces in various poses

with the help of a generic 3D model. We plan to pursue this approach in the near future, and we leave it for the future work.

While this paper was under preparation, the authors became aware of [35], which also proposes a CNN-based approach for face verification/identification and evaluates it on the IJB-A dataset. The method proposed in [35] combines the features from seven independent DCNN models. With finetuning on the JANUS training data and metric learning, our approach works comparable to [35] as shown in Figure 5. Furthermore, with the replacement of ReLU with PReLU and data augmentation, our approach significantly outperforms [35] with only a single model.

### 4.3. Labeled Face in the Wild

We also evaluate our approach on the well-known LFW dataset using the standard protocol which defines 3,000 positive pairs and 3,000 negative pairs in total and further splits them into 10 disjoint subsets for cross validation. Each subset contains 300 positive and 300 negative pairs. It contains 7,701 images of 4,281 subjects. We compare the mean accuracy of the proposed deep model with other state-of-the-

<sup>5</sup>We correct the number reported in [8] previously for the IJB-A identification task because one split of the identification task was performed partially due to the corrupted metadata. (*i.e.* Some images were missing at that time. The current metadata of IJB-A has fixed those errors already.)

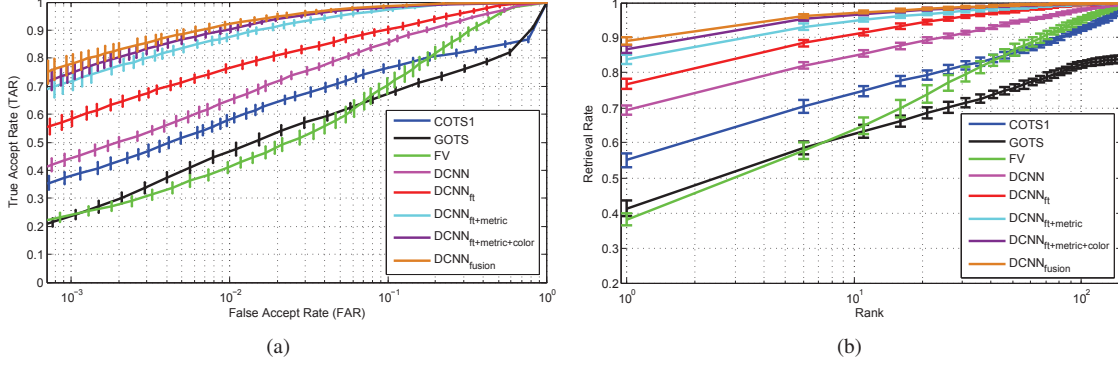


Figure 4. Results on the JANUS CS2 dataset. (a) the average ROC curves and (b) the average CMC curves.

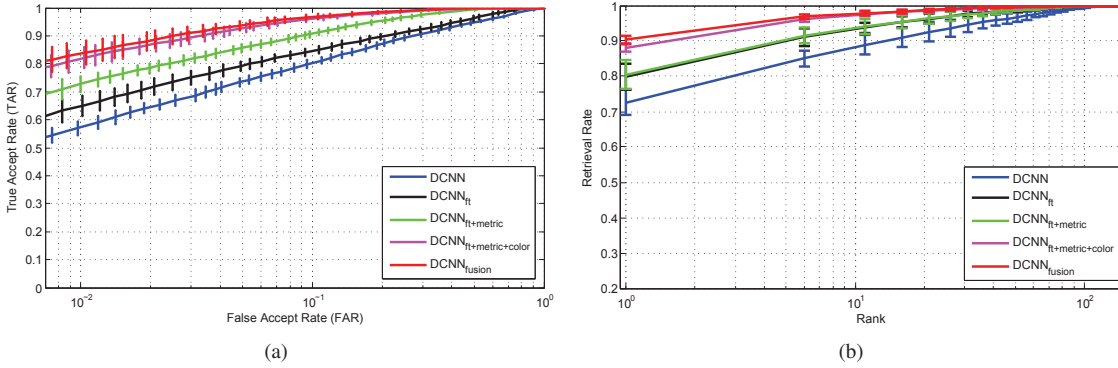


Figure 5. Results on the IJB-A dataset. (a) the average ROC curves for the IJB-A verification protocol and (b) the average CMC curves for IJB-A identification protocol over 10 splits.

| IJB-A-Verif | [35]        | DCNN        | DCNN <sub>ft</sub> | DCNN <sub>ft+m</sub> | DCNN <sub>ft+m+c</sub> | DCNN <sub>fusion</sub> |
|-------------|-------------|-------------|--------------------|----------------------|------------------------|------------------------|
| FAR=1e-2    | 0.732±0.033 | 0.573±0.024 | 0.64±0.045         | 0.787±0.043          | 0.818±0.037            | <b>0.838±0.042</b>     |
| FAR=1e-1    | 0.895±0.013 | 0.8±0.012   | 0.883±0.012        | 0.947±0.011          | 0.961±0.01             | <b>0.967±0.009</b>     |
| IJB-A-Ident | [35]        | DCNN        | DCNN <sub>ft</sub> | DCNN <sub>ft+m</sub> | DCNN <sub>ft+m+c</sub> | DCNN <sub>fusion</sub> |
| Rank-1      | 0.820±0.024 | 0.726±0.034 | 0.799±0.036        | 0.852±0.018          | 0.882±0.01             | <b>0.903±0.012</b>     |
| Rank-5      | 0.929±0.013 | 0.84±0.023  | 0.901±0.025        | 0.937±0.01           | 0.957±0.07             | <b>0.965±0.008</b>     |
| Rank-10     | N/A         | 0.884±0.025 | 0.934±0.016        | 0.954±0.007          | 0.974±0.005            | <b>0.977±0.007</b>     |

Table 3. Results on the IJB-A dataset. The TAR of all the approaches at FAR=0.1 and 0.01 for the ROC curves. The Rank-1, Rank-5, and Rank-10 retrieval accuracies of the CMC curves where subscripts *ft*, *m* and *c* stand for finetuning, metric, and color respectively.

| CS2-Verif | COTS1       | GOTS        | FV[26]      | DCNN        | DCNN <sub>ft</sub> | DCNN <sub>ft+m</sub> | DCNN <sub>ft+m+c</sub> | DCNN <sub>fusion</sub> |
|-----------|-------------|-------------|-------------|-------------|--------------------|----------------------|------------------------|------------------------|
| FAR=1e-2  | 0.581±0.054 | 0.467±0.066 | 0.411±0.081 | 0.649±0.015 | 0.765±0.014        | 0.876±0.013          | 0.904±0.011            | <b>0.921±0.013</b>     |
| FAR=1e-1  | 0.767±0.015 | 0.675±0.015 | 0.704±0.028 | 0.855±0.01  | 0.902±0.011        | 0.973±0.005          | 0.983±0.004            | <b>0.985±0.004</b>     |
| CS2-Ident | COTS1       | GOTS        | FV [26]     | DCNN        | DCNN <sub>ft</sub> | DCNN <sub>ft+m</sub> | DCNN <sub>ft+m+c</sub> | DCNN <sub>fusion</sub> |
| Rank-1    | 0.551±0.03  | 0.413±0.022 | 0.381±0.018 | 0.694±0.012 | 0.768±0.013        | 0.838±0.012          | 0.867±0.01             | <b>0.891±0.01</b>      |
| Rank-5    | 0.694±0.017 | 0.571±0.017 | 0.559±0.021 | 0.809±0.011 | 0.874±0.01         | 0.924±0.009          | 0.949±0.005            | <b>0.957±0.007</b>     |
| Rank-10   | 0.741±0.017 | 0.624±0.018 | 0.637±0.025 | 0.85±0.009  | 0.91±0.008         | 0.949±0.006          | 0.966±0.005            | <b>0.972±0.005</b>     |

Table 4. Results on the JANUS CS2 dataset. The TAR of all the approaches at FAR=0.1 and 0.01 for the ROC curves. The Rank-1, Rank-5, and Rank-10 retrieval accuracies of the CMC curves where subscripts *ft*, *m* and *c* stand for finetuning, metric, and color respectively.

art deep learning-based methods: DeepFace [33], DeepID2 [30], DeepID3 [29], FaceNet [25], Yi *et al.* [38], Wang *et al.* [35], and human performance on the “funneled” LFW images. The results are summarized in Table 5. It can be seen from this table that our approach performs comparably

to other deep learning-based methods. Note that some of the deep learning-based methods compared in Table 5 use millions of data samples for training the model. Whereas we use only the CASIA dataset for training our model which has less than 500K images.

| Method               | #Net | Training Set                                      | Metric                    | Mean Accuracy $\pm$ Std |
|----------------------|------|---|---------------------------|-------------------------|
| DeepFace [33]        | 1    | 4.4 million images of 4,030 subjects, private     | cosine                    | 95.92% $\pm$ 0.29%      |
| DeepFace             | 7    | 4.4 million images of 4,030 subjects, private     | unrestricted, SVM         | 97.35% $\pm$ 0.25%      |
| DeepID2 [30]         | 1    | 202,595 images of 10,117 subjects, private        | unrestricted, Joint-Bayes | 95.43%                  |
| DeepID2              | 25   | 202,595 images of 10,117 subjects, private        | unrestricted, Joint-Bayes | 99.15% $\pm$ 0.15%      |
| DeepID3 [29]         | 50   | 202,595 images of 10,117 subjects, private        | unrestricted, Joint-Bayes | 99.53% $\pm$ 0.10%      |
| FaceNet [25]         | 1    | 260 million images of 8 million subjects, private | L2                        | 99.63% $\pm$ 0.09%      |
| Yi et al. [38]       | 1    | 494,414 images of 10,575 subjects, public         | cosine                    | 96.13% $\pm$ 0.30%      |
| Yi et al.            | 1    | 494,414 images of 10,575 subjects, public         | unrestricted, Joint-Bayes | 97.73% $\pm$ 0.31%      |
| Wang et al. [35]     | 1    | 494,414 images of 10,575 subjects, public         | cosine                    | 96.95% $\pm$ 1.02%      |
| Wang et al.          | 7    | 494,414 images of 10,575 subjects, public         | cosine                    | 97.52% $\pm$ 0.76%      |
| Wang et al.          | 1    | 494,414 images of 10,575 subjects, public         | unrestricted, Joint-Bayes | 97.45% $\pm$ 0.99%      |
| Wang et al.          | 7    | 494,414 images of 10,575 subjects, public         | unrestricted, Joint-Bayes | 98.23% $\pm$ 0.68%      |
| Human, funneled [35] | N/A  | N/A   | N/A                       | 99.20%                  |
| Ours                 | 1    | 490,356 images of 10,548 subjects, public         | cosine                    | 97.15% $\pm$ 0.7%       |
| Ours                 | 1    | 490,356 images of 10,548 subjects, public         | unrestricted, Joint-Bayes | 97.45% $\pm$ 0.7%       |

Table 5. Accuracy of different methods on the LFW dataset.

#### 4.4. Run Time

The DCNN model is trained for about 9 days using NVidia Tesla K40. The feature extraction time takes about 0.006 second per face image. In future, the supervised information will be fed into the intermediate layers to make the model more discriminative and also to converge faster.

#### 5. Conclusion

In this paper, we study the performance of a DCNN method on a newly released challenging face verification dataset, IARPA Benchmark A, which contains faces with full pose, illumination, and other difficult conditions. It was shown that the DCNN approach can learn a robust model from a large dataset characterized by face variations and generalizes well to another dataset. Experimental results demonstrate that the performance of the proposed DCNN on the IJB-A dataset is much better than the FV-based method and other commercial off-the-shelf matchers and is competitive for the LFW dataset.

For future work, we plan to directly train a Siamese network using all the available positive and negative pairs from CASIA-Webface and IJB-A training datasets to fully utilize the discriminative information for realizing better performance.

#### 6. Acknowledgments

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. 2014-14071600012. 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. We thank NVIDIA for donating of the K40 GPU used in this work.

#### References

- [1] 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.
- [2] A. Asthana, S. Zafeiriou, S. Y. Cheng, and M. Pantic. Robust discriminative response map fitting with constrained local models. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3444–3451, 2013.
- [3] A. Asthana, S. Zafeiriou, S. Y. Cheng, and M. Pantic. Incremental face alignment in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1859–1866, 2014.
- [4] X. D. Cao, D. Wipf, F. Wen, G. Q. Duan, and J. Sun. A practical transfer learning algorithm for face verification. In *IEEE International Conference on Computer Vision*, pages 3208–3215. IEEE, 2013.
- [5] D. Chen, X. D. Cao, L. W. Wang, F. Wen, and J. Sun. Bayesian face revisited: A joint formulation. In *European Conference on Computer Vision*, pages 566–579, 2012.
- [6] D. Chen, X. D. Cao, F. Wen, and J. Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2013.
- [7] J.-C. Chen, V. M. Patel, and R. Chellappa. Landmark-based fisher vector representation for video-based face verification. In *IEEE Conference on Image Processing*, 2015.
- [8] J.-C. Chen, V. M. Patel, and R. Chellappa. Unconstrained face verification using deep cnn features. *arXiv preprint arXiv:1508.01722*, 2015.
- [9] J.-C. Chen, S. Sankaranarayanan, V. M. Patel, and R. Chellappa. Unconstrained face verification using fisher vectors



- computed from frontalized faces. In *IEEE International Conference on Biometrics: Theory, Applications and Systems*, 2015.
- [10] Y.-C. Chen, V. M. Patel, R. Chellappa, and P. J. Phillips. Adaptive representations for video-based face recognition across pose. In *IEEE Winter Conference on Applications of Computer Vision*, 2014.
  - [11] Y.-C. Chen, V. M. Patel, P. J. Phillips, and R. Chellappa. Dictionary-based face recognition from video. In *European Conference on Computer Vision*, pages 766–779, 2012.
  - [12] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon. Information-theoretic metric learning. In *International Conference on Machine Learning*, pages 209–216, 2007.
  - [13] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. *arXiv preprint arXiv:1310.1531*, 2013.
  - [14] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, 2014.
  - [15] P. Grother and M. Ngan. Face recognition vendor test(frvt): Performance of face identification algorithms. *NIST Interagency Report 8009*, 2014.
  - [16] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *arXiv preprint arXiv:1502.01852*, 2015.
  - [17] J. Hu, J. Lu, and Y.-P. Tan. Discriminative deep metric learning for face verification in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1875–1882, 2014.
  - [18] Z. Huang, R. Wang, S. Shan, and X. Chen. Projection metric learning on Grassmann manifold with application to video based face recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 140–149, 2015.
  - [19] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACM International Conference on Multimedia*, pages 675–678, 2014.
  - [20] 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 *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
  - [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.
  - [22] J. Lu, V. E. Liong, G. Wang, and P. Moulin. Joint feature learning for face recognition. *IEEE Transactions on Information Forensics and Security*, PP(99):1–1, 2015.
  - [23] O. M. Parkhi, K. Simonyan, A. Vedaldi, and A. Zisserman. A compact and discriminative face track descriptor. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
  - [24] V. M. Patel, T. Wu, S. Biswas, P. J. Phillips, and R. Chellappa. Dictionary-based face recognition under variable lighting and pose. *IEEE Transactions on Information Forensics and Security*, 7(3):954–965, 2012.
  - [25] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. *arXiv preprint arXiv:1503.03832*, 2015.
  - [26] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman. Fisher vector faces in the wild. In *British Machine Vision Conference*, volume 1, page 7, 2013.
  - [27] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
  - [28] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, pages 1988–1996, 2014.
  - [29] Y. Sun, D. Liang, X. Wang, and X. Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015.
  - [30] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. *arXiv preprint arXiv:1412.1265*, 2014.
  - [31] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.
  - [32] Y. Taigman, L. Wolf, and T. Hassner. Multiple one-shots for utilizing class label information. In *British Machine Vision Conference*, pages 1–12, 2009.
  - [33] Y. Taigman, M. Yang, M. A. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014.
  - [34] P. Viola and M. J. Jones. Robust real-time face detection. *International journal of computer vision*, 57(2):137–154, 2004.
  - [35] D. Wang, C. Otto, and A. K. Jain. Face search at scale: 80 million gallery. *arXiv preprint arXiv:1507.07242*, 2015.
  - [36] K. Q. Weinberger, J. Blitzer, and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. In *Advances in neural information processing systems*, pages 1473–1480, 2005.
  - [37] S. Xie, S. G. Shan, X. L. Chen, and J. Chen. Fusing local patterns of gabor magnitude and phase for face recognition. *IEEE Transactions on Image Processing*, 19(5):1349–1361, 2010.
  - [38] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
  - [39] B. C. Zhang, S. G. Shan, X. L. Chen, and W. Gao. Histogram of Gabor phase patterns (hgpp): a novel object representation approach for face recognition. *IEEE Transactions on Image Processing*, 16(1):57–68, 2007.
  - [40] W. Y. Zhao, R. Chellappa, P. J. Phillips, and A. Rosenfeld. Face recognition: A literature survey. *ACM Computing Surveys*, 35(4):399–458, 2003.