# Multi-label Learning Based Deep Transfer Neural Network for Facial Attribute Classification

Ni Zhuang[a], Yan Yan[a,*], Si Chen[b], Hanzi Wang[a], Chunhua Shen[c]

*[a]Fujian Key Laboratory of Sensing and Computing for Smart City,
School of Information Science and Engineering, Xiamen University, Xiamen 361005, China
[b]School of Computer and Information Engineering, Xiamen University of Technology, Xiamen 361024, China
[c]School of Computer Science, The University of Adelaide, Adelaide, SA 5005, Australia*

## Abstract

Deep Neural Network (DNN) has recently achieved outstanding performance in a variety of computer vision tasks, including facial attribute classification. The great success of classifying facial attributes with DNN often relies on a massive amount of labelled data. However, in real-world applications, labelled data are only provided for some commonly used attributes (such as age, gender); whereas, unlabelled data are available for other attributes (such as attraction, hairline). To address the above problem, we propose a novel deep transfer neural network method based on multi-label learning for facial attribute classification, termed FMTNet, which consists of three sub-networks: the Face detection Network (FNet), the Multi-label learning Network (MNet) and the Transfer learning Network (TNet). Firstly, based on the Faster Region-based Convolutional Neural Network (Faster R-CNN), FNet is fine-tuned for face detection. Then, MNet is fine-tuned by FNet to predict multiple attributes with labelled data, where an effective loss weight scheme is developed to explicitly exploit the correlation between facial attributes based on attribute grouping. Finally, based on MNet, TNet is trained by taking advantage of unsupervised domain adaptation for unlabelled facial attribute classification. The three sub-networks are tightly coupled to perform effective facial attribute classification. A distinguishing characteristic of the proposed FMTNet method is that the three sub-networks (FNet, MNet and TNet) are constructed in a similar network structure. Extensive experimental results on challenging face datasets demonstrate the effectiveness of our proposed method compared with several state-of-the-art methods.

*Keywords:* transfer learning, facial attribute classification, multi-label learning, deep learning, convolutional neural networks

## 1. Introduction

Facial attribute classification is an important and fundamental research area in computer vision and pattern recognition. The task of facial attribute classification is to predict the attributes of a facial image, including gender, attraction, race, etc. Recently, facial attribute classification has received increasing attention with a wide range of applications, such as face verification [1, 2, 3], face recognition [4, 5, 6], face retrieval [7]. However, it remains a challenging problem, because of the large facial appearance variations caused by pose, illumination and occlusion, etc.

Early works on facial attribute classification usually characterize the facial attributes based on the histogram representation [2, 3, 8]. For example, Kumar *et al.* [2] propose to firstly extract the low-level features from different regions of a face, and then predict facial attributes with the Support Vector Machine (SVM) for face verification. Cherniavsky *et al.* [8] develop a generative facial feature representation method based on the Haar-like features and investigate a semi-supervised method to predict facial attributes with SVM.

Recent research mainly focuses on using the Deep Neural Network (DNN) to predict facial attributes. Luo *et al.* [9] combine discriminative decision trees with the deep Sum-Product Network (SPN) for facial attribute classification. In [10, 11, 12], the authors firstly extract facial features using DNN and then classify facial attributes with SVM. Ehrlich *et al.* [13] learn the shared feature representation for facial attributes by directly operating on faces and facial landmark points. Rudd *et al.* [14] address the problem of imbalanced data to predict multiple facial attributes.

Generally speaking, methods for facial attribute classification can be divided into two categories: single-label learning based methods [2, 8, 10, 11, 12] and multi-label learning based methods [13, 14]. The single-label learning based methods predict facial attributes separately and thus do not consider the correlation between facial attributes. In contrast, the multi-label learning based methods, which attempt to predict facial attributes simultaneously by using labelled data, have drawn increasing attention. However, in real-world applications, only some commonly used attributes are provided with labelled information, while the other attributes have unlabelled data. Therefore, these methods [13, 14] fail to deal with the facial attribute classification problem when unlabelled information is available (recall that these methods are based on super-

---

*Corresponding author. Tel.:+86-592-2580063
*Email address:* yanyan@xmu.edu.cn (Yan Yan)

vised learning).

Motivated by the above observations, we propose a novel facial attribute classification method, which performs transfer learning based on multi-label learning. More specifically, we take advantage of the transfer DNN technique to predict facial attributes that do not have labelled information in the target domain. To effectively exploit the labelled data in the source domain, we use the multi-label learning technique to predict multiple facial attributes simultaneously, considering the correlation between facial attributes. Fig. 1 shows an illustration of the correlation between different facial attributes. For different learning problems, some carefully designed networks are used, where these networks share the same structure at the former layers of the networks and they only differ at the latter layers. Therefore, the networks can be effectively trained via fine-tuning.

In this paper, we propose an effective deep transfer neural network method, termed FMTNet, which consists of three sub-networks for facial attribute classification. The first sub-network is the Face detection Network (FNet) for face detection. FNet is initialized by using the model learned from a large scale ImageNet dataset [15], and then is fine-tuned by using the facial images. The second sub-network is the Multi-label learning Network (MNet) for facial attribute classification with supervised learning, where multiple attributes are predicted simultaneously. Based on FNet, MNet is fine-tuned by using labelled attributes in the source domain. The network structures at the former layers of both MNet and FNet are the same, whereas the main difference is that multiple fully-connected layers are independently constructed in MNet. The third sub-network is the Transfer learning Network (TNet) for facial attribute classification, when labelled information is not available in the target domain. Based on MNet, TNet makes use of unsupervised domain adaptation to improve the performance of facial attribute classification.

The main contributions of this paper are summarized as follows:

(1) Instead of using single-label learning for each attribute [4, 7], the proposed method effectively performs facial attribute classification based on multi-label learning for the labelled attributes in the source domain. Especially, we propose an effective loss weight scheme to explicitly exploit the correlation between facial attributes based on attribute grouping, which can significantly improve the generalization performance of the proposed method.

(2) Based on multi-label learning, the proposed method leverages transfer learning to predict facial attributes for the unlabelled attributes in the target domain. The transfer neural network successfully transfers the features from the source domain (with labelled information) to the target domain (without labelled information), even when the probability distributions between the two domains are significantly different. Therefore, the proposed method alleviates the dependency on fully labelled training data, especially in the absence of labelled information for some attributes.

The remainder of the paper is organized as follows: In Section 2, some related work is discussed. In Section 3, the details of the proposed FMTNet method for facial attribute classification are described. In Section 4, the experimental results are reported. In Section 5, the conclusions are presented.

## 2. Related Work

The proposed method is closely related to deep learning, multi-label learning and transfer learning. In this section, we briefly discuss the related work.

### 2.1. Deep Learning

The recent great success of deep learning based facial attribute classification is triggered by a growing number of works [10, 11, 12, 16, 17] on learning compact and discriminative features. Compared with the traditional feature extraction methods (such as Scale-Invariant Features Transform (SIFT) [18] and Histogram of Oriented Gradients (HOG) [19]), the deep learning based methods [10, 11, 16] have shown astounding performance improvement. For example, Zhang *et al.* [10] combine the deep features obtained from each poselet of the face region with the deep features of the whole facial image as the final features, and then use SVM for facial attribute classification. Liu *et al.* [11] employ DNN for face localization and apply another neural network to extract features in the face region for attribute classification. Zhong *et al.* [16] extract the features of each attribute based on the hierarchical DNN, and use the linear SVM to classify facial attributes. All the above methods rely on DNN to extract features and require an additional classifier to classify facial attributes. Furthermore, they usually do not consider the correlation between facial attributes. Different from these methods, our proposed method presents an end-to-end network, which effectively exploits the correlation between facial attributes, thus considerably improving the performance of facial attribute classification.

### 2.2. Multi-label Learning

Multi-label learning [20, 21, 22, 23, 24], which aims to learn multiple different but related labels simultaneously, has received much attention so far. For example, Xu *et al.* [46] analyze the local Rademacher complexity of empirical risk minimization (ERM)-based multi-label learning algorithms and then propose a new method that not only results in a sharp generalization error bound, but also provides a tight approximation of the low-rank structure. In [47], the authors present an additional sparse component to deal with the tail labels for the multi-label learning task. In [50], You *et al.* present the privileged multi-label learning (PrML) method to exploit the correlation between labels. For the problem of multi-label learning with missing labels, Jain *et al.* [48] develop a scalable and generative framework, which is based on a latent factor model for the label matrix and an exposure model for missing labels.

Recently, the DNN based multi-label learning method has been proposed for facial attribute classification in [14]. The features learned by DNN exhibit the hierarchical structure (such as pixels, edges, object parts and objects), where the low-level features usually share the similar representation. Therefore, the
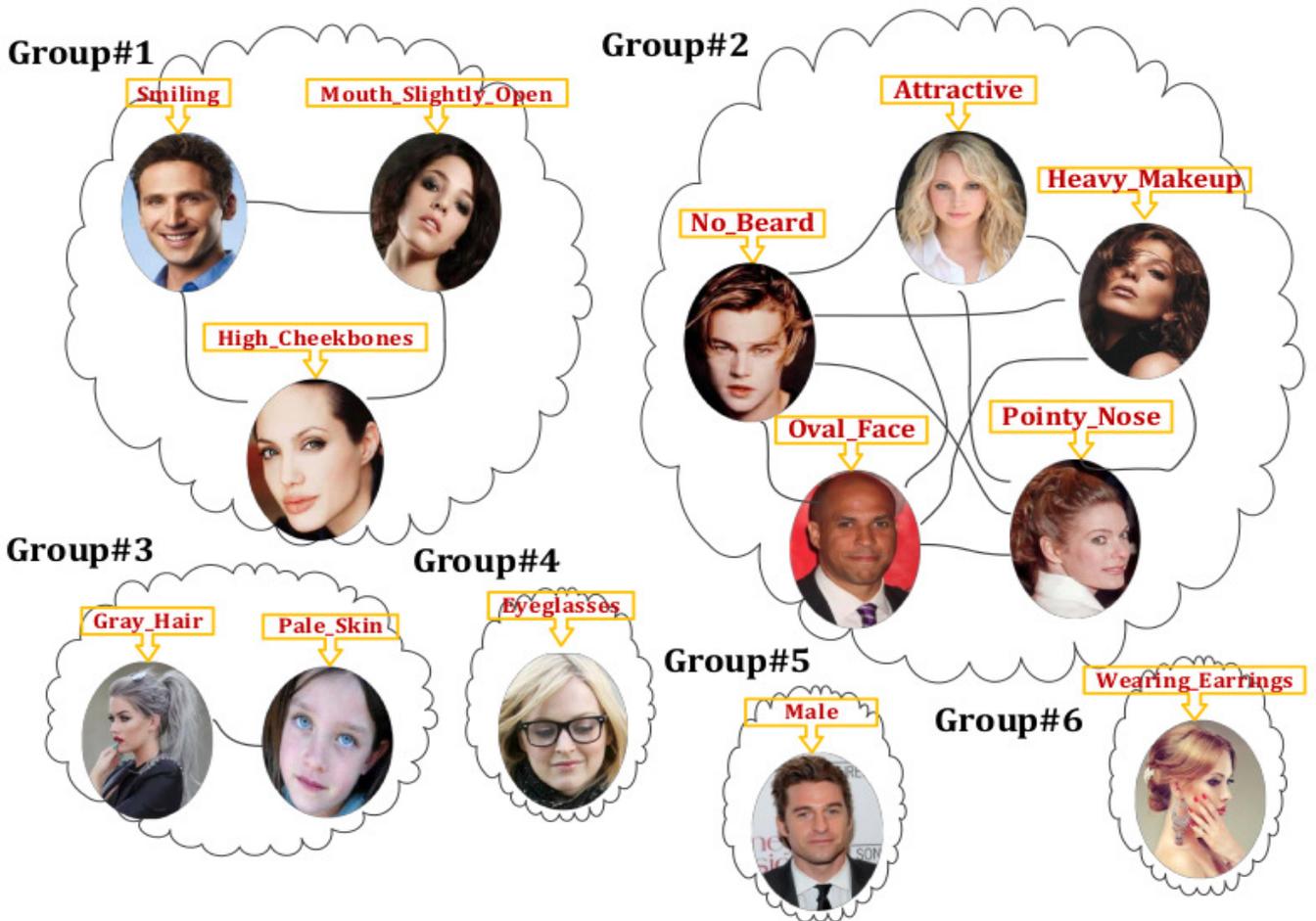
Figure 1: The correlation between facial attributes. Based on the work of Liu *et al.* [11], we obtain the correlation between different attributes. The facial attributes with high correlation can be clustered into the same group.

multi-label learning based on DNN uses the shared features at the low-level layers and separately learns the features of multiple facial attributes at the high-level layers. However, such a method focuses on the imbalance problem of the data and does not explicitly consider the correlation between facial attributes. In this paper, we decompose the multi-label learning problem into a number of binary classification problems, where we solve all these problems using a single neural network at once. Furthermore, we not only consider multiple facial attributes learning, but also investigate the relationship between the different attributes based on attribute grouping.

### 2.3. Transfer Learning

The general performance of the traditional classifier trained with a limited number of labelled data may not be satisfactory, while the manual annotation of abundant training data for various tasks costs too much manpower. Fortunately, transfer learning [25, 26, 27] is an effective technique to improve the performance of the classifier in the target domain given only the annotated data in the source domain, which greatly reduces the labelling cost. Transfer learning also refers to unsupervised domain adaptation [28, 29, 30], and it can adapt the features from the labelled source domain to the unlabelled target domain. Long *et al.* [28] consider both the marginal and conditional distributions between the source and target domains via a Joint Distribution Adaptation (JDA) method. In [29], the authors propose to reduce the domain difference by jointly matching the features and reweighing the instances via a Transfer Joint Matching (TJM) method. Han *et al.* [49] propose a sparse multi-label transfer learning framework, which first learns a multi-label encoded sparse linear embedding space, and then maps the target data onto the learned space for circumventing the problem of multiple tags in practical applications.

With the development of deep learning, the DNN based transfer learning methods [31, 32] usually add the adaptation layers to effectively reduce the discrepancy between the source domain and the target domain. For example, Long *et al.* [32] generalize DNN to domain adaptation, which is formulated as the Deep Adaptation Network (DAN). To the best of our knowledge, our proposed method is the first attempt for the domain adaptation of facial attribute classification. More importantly, we propose to use the deep transfer neural network based on multi-label learning for facial attribute classification.

## 3. The Proposed Method

The proposed facial attribute classification method (i.e., FMTNet) mainly consists of three sub-networks: 1) the Face detection Network (FNet) for face detection; 2) the Multi-label learning Network (MNet) for predicting multiple facial attributes simultaneously; and 3) the Transfer learning Network (TNet) for unlabelled facial attribute classification. These three sub-networks are designed to share the same convolutional layers (i.e., the first 13 convolutional layers in VGG-16 [33]), so that they can be easily fine-tuned. Fig. 2 shows the overall framework of the proposed FMTNet method for facial attribute classification. A detailed description of the proposed method is given in the following subsections.

### 3.1. FNet (Face detection Network)

FNet is designed to perform face detection, which outputs the positions of the faces. In this paper, we mainly follow the work of Faster R-CNN [34] to train FNet. Specifically, as shown in Fig. 2, we first take an image (of any size) as the input and use VGG-16 [33], which has 13 shareable convolutional layers, to obtain the feature map of the whole image. Then, we use the Region Proposal Network (RPN) to obtain a set of rectangular face proposals, each of which has a face score. The RoI pooling layer combines the convolutional features with the predicted bounding boxes, so that FNet can process a theoretical valid back-propagation. Finally, we get the face regions with reference to a series of pre-defined boxes by using the box regression.

The above training scheme shows the discriminative capability of FNet for face detection. Most importantly, FNet is used as the basis of MNet and TNet (i.e., these three sub-networks share the same network structures for their former 13 convolutional layers).

### 3.2. MNet (Multi-label learning Network)

The second sub-network is the Multi-label learning Network (MNet), which simultaneously predicts multiple facial attributes with labelled data. Specifically, we fine-tune the former 13 convolutional layers (which are the same as those in FNet) and train the latter three fully-connected layers for each attribute, as shown in Fig. 2. To reduce the model complexity of MNet, we decrease the dimensions of two fully-connected layers used in VGG-16 [33]. Specifically, we reduce the dimensions of the 'fc14' fully-connected layer from 4,096 to 512, and the dimensions of the 'fc15' fully-connected layer from 4,096 to 256. As shown in the experiments, we still obtain on par results with the lower dimensions of the fully-connected layers, thus reducing computational complexity.

MNet predicts one attribute in conjunction with the other attributes. The structure of MNet (see Fig. 2) contains two parts: (1) the shared layers which are collaboratively trained for all the attributes, and (2) the independent layers which are separately trained for each attribute. In this paper, we tackle the multi-label learning in an attribute-by-attribute style and decompose the multi-label learning problem into a number of binary classification problems. Moreover, in order to improve the generalization performance of MNet, we exploit the correlation between facial attributes based on attribute grouping.

Assume that there are $I$ facial attributes to be predicted. Let us define the softmax loss term, $L_i$, associated with the $i$-th attribute, and the corresponding class label (denoted as $y_i$) is written as follows:

$$L_i = -\frac{1}{N} \sum_{n=1}^{N} log(p_{n,y_i}),  \quad (1)$$

where

$$p_{n,y_i} = \frac{e^{x_{n,y_i}}}{\sum\limits_{y_i=1}^{C_i} e^{x_{n,y_i}}},  \quad (2)$$

and $N$ is the number of training examples; $p_{n,y_i}$ represents the probability of the $i$-th attribute class calculated by the softmax function; $x_{n,y_i}$ represents the value of the last fully-connected layer predicted in the $i$-th attribute class; and $y_i$ can take on $C_i$ values, ranging from 1 to $C_i$.

In this paper, we consider the facial attribute classification problem as a binary classification problem (i.e., $C_i = 2$). In other words, an example, which owns the corresponding label, is treated as a positive sample (i.e., $y_i = 1$); otherwise, it is considered a negative sample (i.e., $y_i = 2$). Of course, each facial attribute classification on MNet is not limited to the binary classification, and it can be easily extended to multi-class classification.

Thus, we derive the loss function of MNet as follow:

$$L = \sum_{i=1}^{I} \lambda_i L_i,  \quad (3)$$

where $\lambda_i$ denotes the loss weight for the $i$-th attribute.

In MNet, multiple attributes are trained together, where each attribute is associated with a loss weight (see Eq. (3)). The loss weight has significant influence on the performance of facial attribute classification (see Section 4.2.1 for the experimental results). Moreover, Liu *et al.* [11] show that the facial attributes have clear grouping patterns. In other words, the facial attributes can be clustered into several groups, where the attributes in one group show high correlation and those between groups have low correlation (see Fig. 1). Therefore, based on this observation, we propose to define the loss weight for each attribute by taking advantage of attribute grouping, which effectively considers the relationship between different attributes. The details of the proposed loss weight scheme are given as follows.

Firstly, we cluster all the $I$ facial attributes into $G$ groups (we use the clustering method in [11]). Suppose that there are $g_m$ attributes (and their corresponding loss weights are defined as $\lambda_{g_1}, \lambda_{g_2}, ..., \lambda_{g_m}$) for the $g$-th group. The loss weights for all the $g_m$ attributes are set to be the same, and their sum is equal to $1/G$. Therefore,

$$\lambda_{g_1} = \lambda_{g_2} = \cdots = \lambda_{g_m} = \frac{1}{G} \cdot \frac{1}{g_m}.  \quad (4)$$
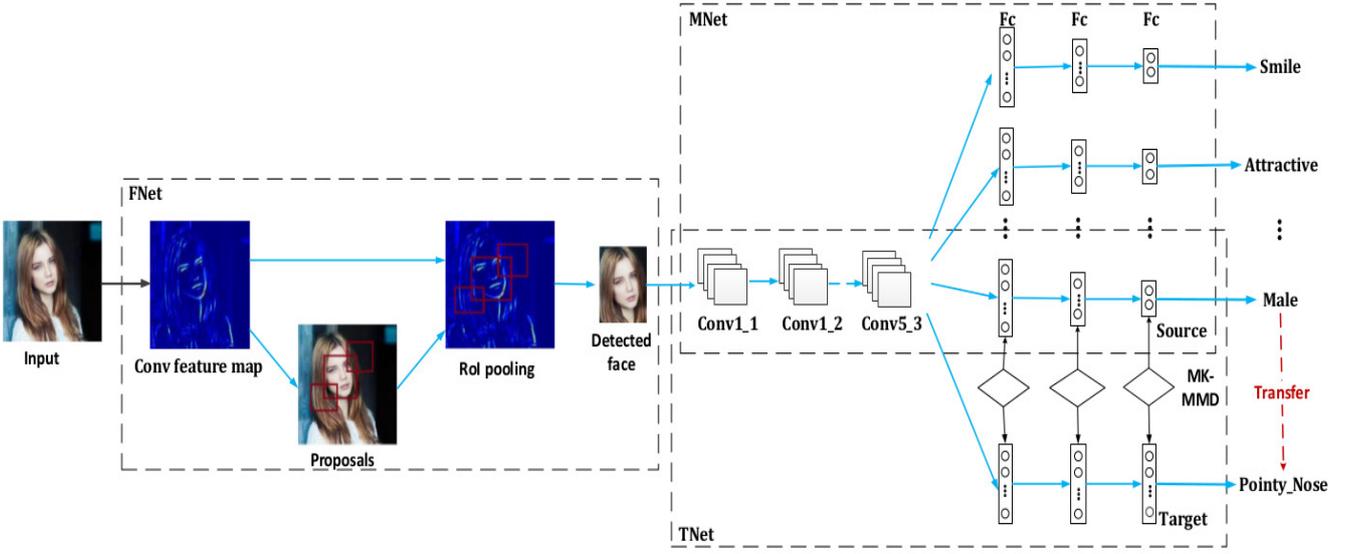
Figure 2: The overall framework of the proposed FMTNet method for facial attribute classification. The first sub-network is FNet for face detection. The second sub-network is MNet for facial attribute classification with supervised learning, where multiple facial attributes are predicted simultaneously. The third sub-network is TNet for facial attribute classification when labelled data are not available. The three sub-networks share the same structure at the former layers but they differ at the latter layers.

If there is only one attribute in a group, we directly set its corresponding loss weight to be *1/G*. For example, we can assume that 6 facial attributes are clustered into 3 groups, where one group (Group 1) has 2 attributes and the other two groups (Group 2 and Group 3) have 3 attributes and 1 attribute, respectively. In this case, the loss weight for each attribute in Group 1 is $1/3 \times 1/2 = 1/6$, while that for each attribute in Group 2 is $1/3 \times 1/3 = 1/9$. The loss weight for the attribute in Group 3 is $1/3$. Thus, the sum of the loss weights for the attributes in each group is the same (i.e., $1/3$). In other words, the proposed loss weight scheme can effectively balance the loss weight of each group (no matter how many attributes are in this group) and prevents TNet from overfitting to one group (especially when many attributes exist in that group). As a result, the performance of MNet can be improved.

Subsequently, we train MNet by using the stochastic gradient descent (SGD) algorithm to minimize the loss function on the training data.

Compared with traditional single-label learning methods [2, 8, 10, 11, 12], which independently train different attributes, MNet effectively overcomes this problem by training multiple attributes simultaneously. Accordingly, the shared layers in MNet significantly reduce the number of network parameters (recall that the DNN based single-label learning methods need to estimate the network parameters for each attribute).

Different from the DNN based multi-label learning method [14], the proposed MNet successfully exploits the correlation between these attributes based on attribute grouping. Recall that we cluster the attributes, which have high correlation, into one group, and then assign the loss weights according to different attribute groups. Therefore, the proposed loss weight scheme takes advantage of attribute grouping to assign the loss weights. It is worth pointing out that attribute grouping is a key step in the proposed method. This is because that not only the proposed loss weight scheme in MNet is based on attribute grouping, but also the performance of TNet is greatly affected by the results of attribute grouping (see Section 4.3.3).

### 3.3. TNet (Transfer learning Network)

The third sub-network is the Transfer learning Network (TNet), which explores the transferability of one attribute with labelled information to another attribute without labelled information. To optimize the training process, TNet shares the similar network structure as MNet and it is fine-tuned based on MNet. In transfer learning, the domain with labelled information is treated as the source domain, and the domain without labelled information is considered as the target domain. Moreover, data in these two domains are usually under different probability distributions. As discussed previously, the attributes in the source domain are trained using the multi-label learning method. In this paper, to deal with different data distributions in the two domains, we introduce a Reproducing Kernel Hilbert Space (RKHS) that is a high-dimensional space, where the domain discrepancy is measured by using the Multi-Kernels Maximum Mean Discrepancies (MK-MMD) criterion proposed by Gretton *et al.* [35]. Fig. 3 gives an intuitive example of transfer learning.

Assume that $X^s = \{x_1^s, x_2^s, ..., x_m^s\}$ consists of $m$ data with the labelled information $Y^s$ in the source domain, and $X^t = \{x_1^t, x_2^t, ..., x_n^t\}$ consists of $n$ data in the target domain. Thus, $D_s = \{(X^s, Y^s)\}$ represents the source domain, and $D_t = \{(X^t)\}$ represents the target domain. Furthermore, let $p$ denote the probability distribution of the source domain and $q$ denote the probability distribution of the target domain. $H_k$ is assumed to be a RKHS defined on a topological space $\chi$ with the characteristic kernel $k$. The MK-MMD $d_k(p, q)$ is defined as the distance
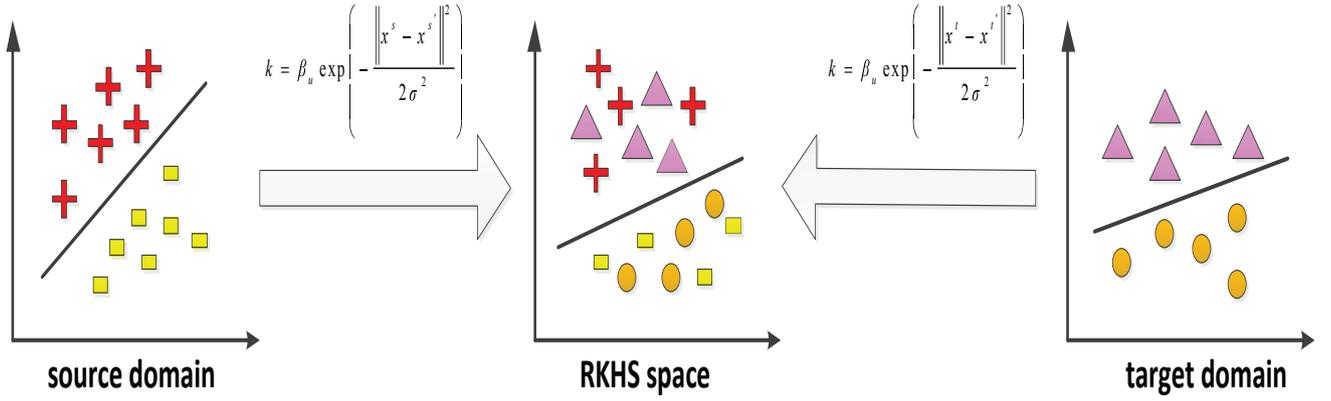
Figure 3: An intuitive example about transfer learning. We use the kernel function to map the samples from the source domain and the target domain into the RKHS space, which can make the two domains closer. In the source domain, the red crosses and the yellow strips denote the positive samples and the negative samples, respectively. In the target domain, the pink triangles and the wheat circles denote the positive samples and the negative samples, respectively.

between the probability distributions $p$ and $q$ in the RKHS. Therefore, we formulate the squared form of MK-MMD as follows:

$$d_k^2(p, q) = \|E_p[\phi(X^s)] - E_q[\phi(X^t)]\|_{H_k}^2, \quad (5)$$

where $E_p[\cdot]$ is the mean of $p$ and $E_q[\cdot]$ is the mean of $q$. $\phi(\cdot)$ is a feature mapping function, which maps the features from the original feature space to the RKHS space [35].

In MK-MMD, we select a kernel from a particular family of kernels $K$ (here, the multiple Gaussian kernels are used for optimal kernel selection). Let us denote $\{k_u\}_{u=1}^d$ as a set of positive definite functions, where $k_u: \chi \times \chi \to R$. Hence,

$$K = \{k : k = \sum_{u=1}^d \beta_u k_u, \sum_{u=1}^d \beta_u = 1, \beta_u \geqslant 0, \forall u \in \{1, ..., d\}\}, \quad (6)$$

where $\beta_u$ is the coefficient of the $u$-th kernel. We constrain $\beta_u$ to guarantee that the derived kernel is characteristic.

The characteristic kernel, $k(\cdot)$, is defined as $k(x^s, x^t) = <\phi(x^s), \phi(x^t)>$. Accordingly, the distance between two probability distributions is computed as follows:

$$d_k^2(p, q) = E_{x^s, x^{s'}}[k(x^s, x^{s'})] - 2E_{x^s, x^t}[k(x^s, x^t)]$$
$$+ E_{x^t, x^{t'}}[k(x^t, x^{t'})], \quad (7)$$

where $x^s$, $x^{s'}$ are two features in the source domain and $x^t$, $x^{t'}$ are two features in the target domain. $E_{x^s, x^{s'}}[\cdot]$ is the mean of $x^s$ and $x^{s'}$, and $E_{x^s, x^t}[\cdot]$ and $E_{x^t, x^{t'}}[\cdot]$ are similarly defined.

The objective of transfer learning is to minimize the domain discrepancy between the source domain and the target domain, which is effectively measured by the distance between the probability distributions from the source domain and the target domain (i.e., MK-MMD). Therefore, we have:

$$\min_{D_s, D_t} D_F(D_s, D_t) = \min_{p, q} d_k^2(p, q), \quad (8)$$

where $D_F (= d_k^2(p, q))$ denotes the domain discrepancy between two domains in the fully-connected layer.

We respectively embed three MK-MMD loss layers in the

last three layers of TNet, and thus the total loss of the whole TNet is the sum of the softmax loss and three MK-MMD losses, which can be written as:

$$L = \sum_{i=1}^3 D_{F_i} + \alpha L_s, \quad (9)$$

where $L$ is the total loss of the whole TNet, which is trained by the SGD method. The $D_{F_i}$ denotes the MK-MMD loss in the last $i$-th layer, and $L_s$ represents the softmax loss of the source domain ($\alpha$ is the loss weight).

The loss of TNet (see Eq. (9)) consists of two terms (i.e., the MK-MMD loss and the softmax loss). The objective of minimizing the MK-MMD loss is to reduce the discrepancy between the distributions of the source and target domains, while the objective of minimizing the softmax loss is to enhance the discriminability of the features in the source domain. In this paper, the MK-MMD loss is based on the kernel mean matching, where only features of data are used (while the labeled information is not used). Meanwhile, the softmax loss is defined in the source domain, where the labelled information is considered. By combining these two terms (i.e., the MK-MMD loss and the softmax loss), the discriminability of the features in the target domain is enhanced.

An example is given in Fig. 2, where the features learned from the 'Male' attribute are transferred to the 'Pointy_Nose' attribute. The 'Male' attribute has a large amount of labelled training data in the source domain; whereas, the 'Pointy_Nose' attribute has unlabelled training data in the target domain. Although the network structures of these two attributes are the same, the discrepancy of the features obtained in the latter layers is large. As shown in the experiments, the performance of feature transferability significantly drops if we perform direct transfer (i.e., train the network with the source data having labelled information, and then directly extract the features of the target data for classification). However, based on the MK-MMD loss function, the domain discrepancy between the 'Male' attribute and the 'Pointy_Nose' attribute is effectively measured. As a result, by minimizing the MK-MMD loss, the

difference between the two attributes is greatly reduced so that the feature transferability is significantly improved.

It is worth pointing out that TNet and MNet share the same feature extraction layers (i.e., the former 13 convolutional layers), while the latter fully-connected layers are different (see Fig. 2 for an illustration). Considering the limited number of training data, TNet directly uses the parameters of the first 6 convolutional layers, while fine-tuning the parameters of the subsequent 7 convolutional layers of MNet (i.e., we freeze 'conv1_1' - 'conv3_2' to learn generic features and fine-tune 'conv3_3' - 'conv5_3' to correct the slight domain biases on the convolutional layers of VGG-16 [33]). This is because that the fine-tuning strategy is usually much easier to obtain the network parameters which are learned by taking advantage of multi-label learning based on the labelled facial attributes data.

Although recent works [31, 32] learn the transferable features with DNN, which also uses the MK-MMD loss function, the proposed TNet has some differences: 1) The proposed TNet focuses on facial attribute classification, while the methods in [31, 32] are proposed for office image classification. 2) The proposed TNet is fine-tuned based on MNet and thus the network structures of TNet and MNet are similar. As a result, TNet takes both advantages of transfer learning and multi-label learning. In contrast, the methods in [31, 32] mainly perform transfer learning via the adaptation network.

### 3.4. The Overall FMTNet

In the previous subsections, we have developed all the ingredients for the proposed FMTNet method. Specifically, we use FNet for face detection, use MNet for predicting multiple facial attributes simultaneously with labelled data, and use TNet for predicting facial attributes without labelled information. The training process of the overall FMTNet method is described in algorithm 1.

## 4. Experiments

In this section, we show the effectiveness of the proposed method for facial attribute classification. In Section 4.1, the datasets and parameter settings are described. In Section 4.2, the performance of MNet is evaluated. The influence of loss weight is given in Section 4.2.1. The comparison between single-label learning and multi-label learning is shown in Section 4.2.2. Performance comparison with several state-of-the-art multi-label learning methods is described in Section 4.2.3. Finally, in Section 4.3, the performance of the proposed FMTNet is evaluated. The performances obtained by direct transfer and transferring under different correlations are given in Sections 4.3.1 and 4.3.2, respectively. Performance comparison with several state-of-the-art transfer learning methods is described in Section 4.3.3.

### 4.1. Datasets and Parameter Settings

We train and evaluate different facial attribute classification methods on the CelebA [36] and LFWA [37] datasets. The CelebA dataset has 202,599 facial images of 10,177 identities,

---

**Algorithm 1** The training process of the proposed FMTNet method

---

**Input:** The training datasets $\mathcal{D}_f$, $\mathcal{D}_m$, and $\mathcal{D}_t$ for FNet, MNet and TNet, respectively.
**Output:** The network parameters $\mathcal{W}_f$, $\mathcal{W}_m$, and $\mathcal{W}_t$ for FNet, MNet and TNet, respectively.
1: **Initialization:** Initialize $\mathcal{W}_f$ with VGG-16;
2: **while** not converge **do**
3:     **for** each training sample $\mathbf{x}_i \in \mathcal{D}_f$ **do**
4:         Forward pass to obtain the feature representation of $\mathbf{x}_i$;
5:         Back propagate to update the network parameters $\mathcal{W}_f$;
6:     **end for**
7: **end while**
8: Cluster all the attributes into $G$ groups;
9: Calculate the loss weight for each attribute via Eq. (4);
10: **Initialization:** Initialize $\mathcal{W}_m$ with $\mathcal{W}_f$;
11: **while** not converge **do**
12:     **for** each training sample $\mathbf{x}_i \in \mathcal{D}_m$ **do**
13:         Forward pass to obtain the feature representation of $\mathbf{x}_i$;
14:         Back propagate to update the network parameters $\mathcal{W}_m$ via Eq. (3);
15:     **end for**
16: **end while**
17: **Initialization:** Initialize $\mathcal{W}_t$ with $\mathcal{W}_m$;
18: **while** not converge **do**
19:     **for** each training sample $\mathbf{x}_i \in \mathcal{D}_t$ **do**
20:         Forward pass to obtain the feature representation of $\mathbf{x}_i$;
21:         Back propagate to update the network parameters $\mathcal{W}_t$ via Eq. (9);
22:     **end for**
23: **end while**

---

and 40 binary attribute annotations are provided for each facial image. The CelebA dataset is divided into three parts: training, validation and test. More specifically, the training set, the validation set and the test set respectively contain 162,770 images, 19,867 images and 19,962 images. The LFWA dataset, with more than 1,680 identities, contains more than 13,000 facial images collected from the web. Each image in this dataset has 73 binary attribute annotations, where the positive or negative values indicate the presence or absence of the corresponding attributes, respectively. We use half of the LFWA dataset for training and half of the LFWA dataset for test.

Recall that the proposed method focuses on facial attribute classification, where the detected faces are used to perform subsequent classification. Therefore, we do not compare FNet with other face detection methods in this paper. In fact, the performance of the Faster R-CNN for face detection is reported in [38]. In this experiment, we use the training data from the CelebA dataset to train FNet and then use FNet as the basis of both MNet and TNet.

Table 1: The classification accuracy (%) obtained by MNet with different loss weight schemes. 'MNet_Eq', 'MNet_Emp' and 'MNet_Prop' represent MNet with the same loss weight scheme, MNet with the emphasized loss weight scheme and MNet with the proposed loss weight scheme, respectively. The best results are in boldface.

| Attributes | CelebA | | | LFWA | | |
|---|---|---|---|---|---|---|
| | MNet_Eq | MNet_Emp | MNet_Prop | MNet_Eq | MNet_Emp | MNet_Prop |
| High_Cheekbone | 87.03 | 87.04 | **88.19** | 84.29 | 82.34 | **85.79** |
| Mouth_Open | 93.72 | 93.59 | **94.16** | 77.92 | 75.27 | **81.59** |
| Smiling | 92.39 | 92.51 | **93.21** | 88.57 | 83.68 | **89.49** |
| Attractive | 82.46 | 83.09 | **83.29** | 73.93 | 77.59 | **77.78** |
| Bangs | 95.94 | **96.08** | 96.03 | 87.62 | **89.90** | 89.42 |
| Blond_Hair | 95.73 | **96.20** | 96.08 | 96.85 | 96.65 | **97.04** |
| Brown_Hair | 88.39 | 89.40 | **89.58** | 73.45 | 76.59 | **78.54** |
| Heavy_Makeup | 91.21 | **91.90** | 91.87 | 92.86 | 93.87 | **93.99** |
| No_Beard | 96.19 | 96.46 | **96.52** | 80.00 | 77.99 | **80.52** |
| Oval_Face | 74.64 | **76.57** | 76.38 | 71.61 | 71.96 | **73.69** |
| Pointy_Nose | 75.63 | 77.73 | **77.79** | 76.49 | 79.85 | **82.13** |
| Rosy_Cheeks | 94.80 | **95.42** | 95.37 | 81.55 | 85.49 | **85.51** |
| Wavy_Hair | 84.55 | **85.91** | 85.62 | 76.25 | 78.48 | **80.24** |
| Lipstick | 93.97 | 94.10 | **94.20** | 92.68 | 92.71 | **93.20** |
| Young | 87.66 | **89.18** | 88.70 | 85.42 | 85.68 | **85.98** |
| Gray_Hair | 98.03 | 98.28 | **98.33** | 88.51 | 88.44 | **90.21** |
| Pale_Skin | 96.68 | 96.97 | **97.17** | 83.87 | 69.64 | **90.86** |
| Blurry | 96.11 | 96.25 | **96.42** | 83.87 | 83.35 | **85.18** |
| Black_Hair | 89.14 | 89.92 | **90.29** | 90.00 | 86.22 | **90.49** |
| Straight_Hair | 83.67 | 84.51 | **84.92** | 71.61 | 68.54 | **77.26** |
| Eyeglasses | 99.68 | 99.62 | **99.71** | 91.37 | 88.33 | **92.32** |
| Hat | 99.06 | 99.09 | **99.18** | 90.60 | 85.98 | **90.79** |
| 5 O.C. Shadow | 94.31 | 94.72 | **94.98** | 74.17 | 70.10 | **75.51** |
| Bald | 98.89 | 98.91 | **99.01** | 91.07 | 89.14 | **92.51** |
| Goatee | 97.13 | 97.41 | **97.63** | 79.94 | 76.89 | **81.40** |
| Male | 97.65 | 98.11 | **98.48** | **92.44** | 90.45 | 92.20 |
| Mustache | 96.92 | 96.91 | **97.03** | 90.48 | 87.05 | **92.40** |
| Sideburns | 97.53 | 97.86 | **97.96** | 77.38 | 73.36 | **79.46** |
| Necktie | 97.05 | 97.06 | **97.08** | 80.30 | 78.91 | **80.79** |
| Arched_Eyebrow | 81.66 | 83.15 | **83.89** | **78.21** | 75.68 | 78.18 |
| Bags_Under_Eye | 84.60 | 84.85 | **85.46** | 78.45 | 76.07 | **79.30** |
| Big_Lips | 70.06 | 71.36 | **71.76** | 72.68 | 68.20 | **74.48** |
| Big_Nose | 83.78 | 84.41 | **84.56** | 80.06 | 77.05 | **82.02** |
| Bushy_Eyebrow | 92.20 | 92.29 | **92.80** | 70.30 | 68.47 | **75.24** |
| Chubby | 95.22 | **95.81** | 95.76 | 73.93 | 73.26 | **74.73** |
| Double_Chin | 95.89 | **96.42** | **96.42** | **79.46** | 75.71 | 79.17 |
| Narrow_Eyes | 85.75 | 86.94 | **87.48** | 77.26 | 72.35 | **78.05** |
| Recede_Hair | 93.70 | 93.72 | **93.82** | 83.75 | 82.78 | **84.11** |
| Earring | 90.14 | 90.46 | **91.04** | 92.74 | 91.96 | **93.17** |
| Necklace | 86.69 | 87.71 | **88.15** | 88.15 | 87.90 | **88.81** |
| Average | 90.90 | 91.45 | **91.66** | 82.50 | 80.85 | **84.34** |

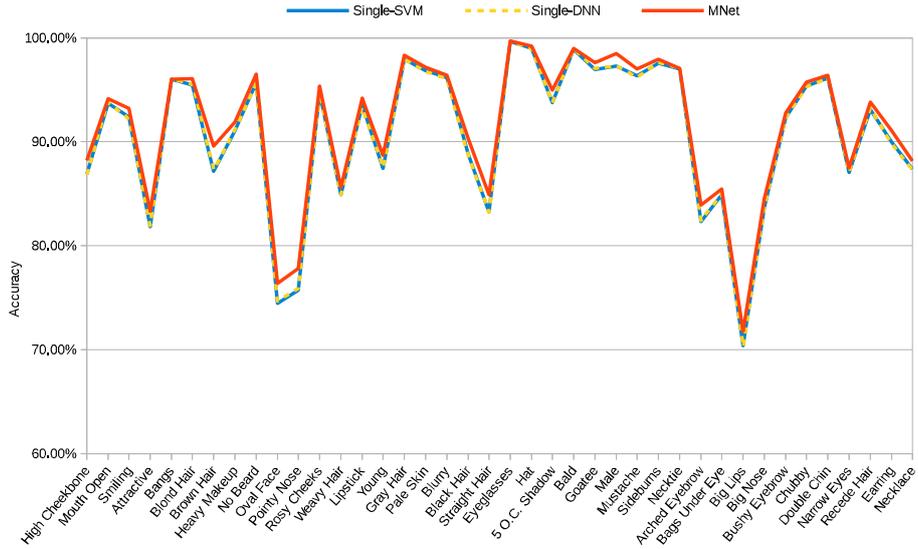### 4.2. Results on the Multi-label Learning Network

In this experiment, we evaluate the performance of MNet using 40 attributes, as done in [11]. The detailed setting of MNet is given as follows: Firstly, MNet is initialized based on FNet. Then, we use 162,770 images on the CelebA dataset to train MNet. Finally, we test MNet with 19,962 images on the CelebA dataset. For the LFWA dataset, we use 6,571 images to train MNet and use 6,571 images to evaluate its performance.
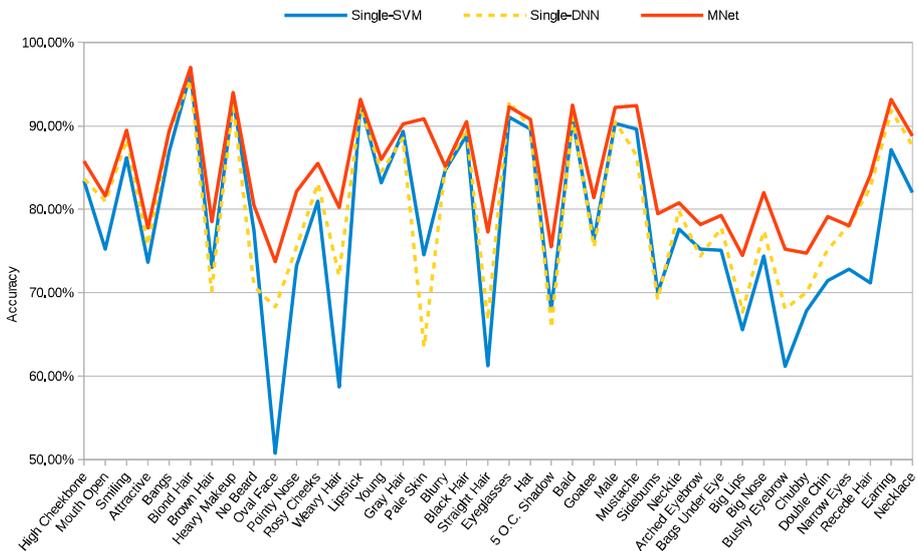
#### 4.2.1. Influence of Loss Weight

In this subsection, we evaluate the performance of MNet with different loss weight schemes on the CelebA and LFWA datasets, respectively. In Table 1, we show the classification accuracy obtained by MNet with different loss weight schemes, where MNet with the same loss weight scheme (i.e., the loss

weight for each attribute is set to 1.0) and MNet with the proposed loss weight scheme are represented by 'MNet_Eq' and 'MNet_Prop', respectively. In 'MNet_Emp', we emphasize the attributes of one group among all groups (i.e., the attributes in one group are assigned with higher loss weights than those in the other groups). In particular, in this experiment we choose one group containing multiple attributes (e.g., the group 2 shown in Fig. 1) and assign the large loss weight (i.e., 1.0) to each attribute in this group. We assign the small loss weight (i.e., 0.1) to each attribute in the other different groups.

From the experimental results in Table 1, we can see that 'MNet_Prop' achieves better results than 'MNet_Eq' on both CelebA and LFWA datasets. In comparison to 'MNet_Eq', the improvement of 'MNet_Prop' on the LFWA dataset is more obvious (about 2% improvement). The improvement is

(a) The CelebA dataset



(b) The LFWA dataset

Figure 4: Performance comparison between the two methods based on single-label learning and MNet based on multi-label learning on (a) the CelebA dataset and (b) the LFWA dataset.

more significant for some individual attributes (for example, 'MNet_Prop' achieves about 4% higher classification accuracy for the 'Attractive' attribute with regard to 'MNet_Eq' on the LFWA dataset ). The proposed 'MNet_Prop' effectively exploits the correlation between facial attributes by taking advantage of attribute grouping. More specifically, the loss weights corresponding to different facial attributes are assigned according to the results of attribute grouping (see Eq. (4)). In this manner, 'MNet_Prop' balances the loss weights for different facial attributes, which not only considers the differences among the facial attributes (note that 'MNet_Eq' treats each facial attribute equally), but also avoids over-fitting to one group (note that 'MNet_Emp' tends to focus on training one group, leading to the problem of over-fitting).

Note that the performance of facial attribute classification is greatly affected by large facial appearance variations caused by pose, illumination, occlusion, etc. However, some facial attributes are more sensitive to the facial appearance variations than the other facial attributes. In other words, the classification task for some facial attributes is more difficult than that for the other facial attributes. For example, the variations caused by the face pose have more influence on the 'Pointy_Nose' attribute than the 'Gray_Hair' attribute. Thus, the classification performance obtained by the proposed method on the 'Gray_Hair' attribute (about 98% classification accuracy on the CelebA dataset) is much better than that on the 'Pointy_Nose' attribute (about 78% classification accuracy on the CelebA dataset).

9

#### 4.2.2. Single-label Learning vs. Multi-label Learning

Secondly, we compare the proposed MNet with two methods (Single-SVM and Single-DNN) that are based on single-label learning (each attribute is learnt separately) on the CelebA and LFWA datasets, respectively. Specifically, Single-SVM denotes the method that uses DNN to learn the features of facial attributes, and then applies an SVM classifier for classification. 'Single-DNN' represents the method that uses DNN to extract facial features and predict facial attributes. The classification accuracy comparison is reported in Fig. 4.

From Fig. 4, we can observe that the proposed MNet achieves better results than the two competing methods that are based on single-label learning, which validates the effectiveness of the proposed method. This is because that multi-label learning effectively exploits the correlation between different facial attributes. In contrast, the network based on single-label learning does not consider the relationship among the facial attributes. Therefore, the single-label learning based methods do not make full use of the intrinsic information of facial attributes.

In summary, compared with the single-label learning, the multiple attributes can be simultaneously exploited in MNet. Furthermore, the proposed loss weight scheme assigns different loss weights for different attributes, which makes MNet based on multi-label learning more effective than the network based on single-label learning.

#### 4.2.3. Comparison with the State-of-the-art Methods

Finally, we compare the performance of MNet with the following state-of-the-art facial attribute classification methods trained with the labelled information: FaceTracer [39], two versions of PANDA [10] (i.e., PANDA-w and PANDA-l), two versions of LNets+ANet [11] (i.e., LNets+ANet (w/o) (without pre-training) and LNets+ANet), and MT-RBM (PCA) [13]. FaceTracer [39] uses the features of HOG and the color histogram to train an SVM for facial attribute classification. PANDA [10] employs hundreds of poselets, which are aligned to predict facial attributes. LNets + ANet [11] cascades two deep neural networks to detect the face and then learns the facial attributes from the detected parts. MT-RBM (PCA) [13] learns the joint feature representation of the faces and facial landmark points for predicting facial attributes.

Classification accuracies obtained by all the competing methods are reported in Table 2. For the CelebA dataset, we compare MNet with all the competing methods. The results obtained by the competing methods on the CelebA dataset are taken from [13]. For the LFWA dataset, we compare MNet with all the competing methods except for MT-RBM (PCA) [13], because that MT-RBM (PCA) [13] does not release the code on the LFWA dataset. Moreover, the LFWA dataset is not suitable for MT-RBM (PCA) to train and test, since the scale of the LFWA dataset is small and the data distribution of the LFWA dataset is unbalanced. The results obtained by the competing methods on the LFWA dataset are taken from [11]. We follow the same evaluation protocol provided in [11, 13].

As shown in Table 2, MNet significantly outperforms the other competing methods on the CelebA dataset. Moreover, MNet achieves similar average accuracy compared with LNet + ANet (and obtains much better results than FaceTracer, PANDA-w and PANDA-l) on the LFWA dataset. LNet + ANet uses multiple patches cropped from the face region for data augmentation, which generates much more training data than our proposed method (using half of the images on the LFWA dataset). Thus, the training complexity of LNet + ANet [11] is much higher than that of MNet. In general, MNet achieves better or comparable performance compared with the state-of-the-art methods.

### 4.3. Results of the Proposed Method

In this section, we evaluate the performance of the proposed FMTNet method for transfer learning on the CelebA dataset. The reason why we do not evaluate the proposed FMTNet on the LFWA dataset is that the data distribution of this dataset is unbalanced (for example, almost all the images are labelled with the 'Male' attribute and the 'Pointy_Nose' attribute, while only a few images are labelled with the 'Heavy_Makeup' attribute), making the LFWA dataset difficult to be used for evaluating the performance of transfer learning.

#### 4.3.1. Direct Transfer

Firstly, we perform the direct transfer method from one attribute to other attributes on the CelebA dataset. In other words, we train a model with the labelled data of one facial attribute in the source domain, and then we use the trained model to directly predict other facial attributes in the target domain. The classification accuracy obtained on one attribute, which is directly transferred to other attributes on the CelebA dataset is reported in Fig. 5. Specifically, we use 162,770 images to train the 'Attractive' attribute which has the labelled information (while other attributes can also be used, and we observe similar results). Then, we predict the other attributes with this trained model.

As shown in Fig. 5, because the model is trained with the labelled information of the 'Attractive' attribute, the best results are obtained in predicting the 'Attractive' attribute. Furthermore, the accuracy obtained on the three attributes – 'Young', 'Lipstick' and 'Heavy_Makeup' – is slightly better than that obtained on the other attributes, because these three attributes belong to the same group as the 'Attractive' attribute. However, the results for other transfer tasks are worse, since we train the model without using the labelled information of these attributes. It is worth pointing out that although the 'No_Beard', 'Oval_Face' and 'Pointy_Nose' attributes belong to the same group as the 'Attractive' attribute, the classification accuracy obtained on these attributes is not high. Therefore, the direct transfer method is not desirable and the performance is not stable.

#### 4.3.2. Transfer under Different Correlations

The correlation between the source and target domains will greatly affect the performance of unlabelled facial attribute classification. In this subsection, we evaluate the performance of FMTNet under different correlations between the source and target domains.

Table 2: The classification accuracy (%) comparison between MNet and several state-of-the-art methods for 40 facial attributes on the CelebA dataset and the LFWA dataset. The best results are in boldface.

| Attributes | CelebA | | | | | | | | LFWA | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FaceTracer [39] | PANDA-w [10] | PANDA-l [10] | ANet [40] | LNets+ANet(w/0) [11] | LNets+ANet [11] | MT-RBM (PCA) [13] | MNet | FaceTracer [39] | PANDA-w [10] | PANDA-l [10] | ANet [40] | LNets+ANet(w/0) [11] | LNets+ANet [11] | MNet |
| High_Cheekbone | 84 | 80 | 86 | 85 | 84 | 87 | 83 | **88.19** | 77 | 75 | 86 | 79 | 83 | **88** | 85.79 |
| Mouth_Open | 87 | 82 | 93 | 85 | 86 | 92 | 82 | **94.16** | 77 | 74 | 78 | 76 | 78 | **82** | 81.59 |
| Smiling | 89 | 89 | 92 | 92 | 88 | 92 | 88 | **93.21** | 78 | 77 | 89 | 82 | 88 | **91** | 89.49 |
| Attractive | 78 | 77 | 81 | 79 | 77 | 81 | 76 | **83.29** | 71 | 70 | 81 | 75 | 80 | **83** | 77.78 |
| Bangs | 88 | 89 | 92 | 94 | 92 | 95 | 88 | **96.03** | 72 | 79 | 84 | 84 | 84 | 88 | **89.42** |
| Blond_Hair | 80 | 81 | 93 | 86 | 91 | 95 | 91 | **96.08** | 88 | 87 | 94 | 90 | 94 | 97 | **97.04** |
| Brown_Hair | 60 | 69 | 77 | 74 | 78 | 80 | 83 | **89.58** | 62 | 65 | 74 | 71 | 73 | 77 | **78.54** |
| Heavy_Makeup | 85 | 84 | 90 | 87 | 85 | 90 | 85 | **91.87** | 88 | 86 | 93 | 89 | 91 | **95** | 93.99 |
| No_Beard | 90 | 87 | 93 | 91 | 92 | 95 | 90 | **96.52** | 69 | 63 | 75 | 69 | 75 | 79 | **80.52** |
| Oval_Face | 64 | 62 | 65 | 65 | 63 | 66 | 73 | **76.38** | 66 | 64 | 72 | 66 | 71 | **74** | 73.69 |
| Pointy_Nose | 68 | 65 | 71 | 67 | 70 | 72 | 73 | **77.79** | 74 | 68 | 76 | 72 | 76 | 80 | **82.13** |
| Rosy_Cheeks | 84 | 81 | 87 | 85 | 87 | 90 | 94 | **95.37** | 70 | 64 | 73 | 71 | 72 | 78 | **85.51** |
| Wavy_Hair | 73 | 76 | 77 | 79 | 75 | 80 | 72 | **85.62** | 62 | 63 | 75 | 65 | 73 | 76 | **80.24** |
| Lipstick | 89 | 88 | 93 | 91 | 90 | 93 | 89 | **94.20** | 87 | 83 | 93 | 86 | 92 | **95** | 93.20 |
| Young | 80 | 77 | 84 | 81 | 83 | 87 | 81 | **88.70** | 80 | 76 | 82 | 79 | 82 | **86** | 85.98 |
| Gray_Hair | 90 | 88 | 94 | 93 | 93 | 97 | 97 | **98.33** | 78 | 77 | 81 | 82 | 81 | 84 | **90.21** |
| Pale_Skin | 83 | 84 | 91 | 89 | 87 | 91 | 96 | **97.17** | 70 | 64 | 84 | 68 | 81 | 84 | **90.86** |
| Blurry | 81 | 77 | 86 | 83 | 80 | 84 | 95 | **96.42** | 73 | 70 | 74 | 75 | 70 | 74 | **85.18** |
| Black_Hair | 70 | 74 | 85 | 77 | 84 | 88 | 76 | **90.29** | 76 | 78 | 87 | 82 | 86 | 90 | **90.49** |
| Straight_Hair | 63 | 67 | 69 | 70 | 69 | 73 | 80 | **84.92** | 67 | 68 | 73 | 72 | 71 | 76 | **77.26** |
| Eyeglasses | 98 | 94 | 98 | 96 | 96 | 99 | 96 | **99.71** | 90 | 84 | 89 | 88 | 92 | **95** | 92.32 |
| Hat | 89 | 91 | 96 | 93 | 96 | 99 | 97 | **99.18** | 75 | 78 | 82 | 82 | 84 | 88 | **90.79** |
| 5 O.C. Shadow | 85 | 82 | 88 | 86 | 88 | 91 | 90 | **94.98** | 70 | 64 | 84 | 78 | 81 | **84** | 75.51 |
| Bald | 89 | 92 | 96 | 92 | 95 | 98 | 98 | **99.01** | 77 | 82 | 84 | 86 | 83 | 88 | **92.51** |
| Goatee | 93 | 86 | 93 | 92 | 92 | 95 | 96 | **97.63** | 69 | 65 | 75 | 68 | 75 | 78 | **81.40** |
| Male | 91 | 93 | 97 | 95 | 94 | 98 | 90 | **98.48** | 84 | 86 | 92 | 91 | 91 | **94** | 92.20 |
| Mustache | 91 | 83 | 93 | 87 | 91 | 95 | 97 | **97.03** | 83 | 77 | 87 | 79 | 87 | 92 | **92.40** |
| Sideburns | 94 | 90 | 93 | 94 | 91 | 96 | 96 | **97.96** | 71 | 68 | 76 | 72 | 72 | 77 | **79.46** |
| Necktie | 86 | 88 | 91 | 90 | 86 | 93 | 94 | **97.08** | 71 | 70 | 79 | 72 | 76 | 79 | **80.79** |
| Arched_Eyebrow | 76 | 73 | 78 | 75 | 74 | 79 | 77 | **83.89** | 67 | 63 | 79 | 66 | 78 | **82** | 78.18 |
| Bags_Under_Eye | 76 | 71 | 79 | 77 | 73 | 79 | 81 | **85.46** | 65 | 63 | 80 | 72 | 79 | **83** | 79.30 |
| Big_Lips | 64 | 61 | 67 | 63 | 66 | 68 | 69 | **71.76** | 68 | 64 | 73 | 70 | 72 | **75** | 74.48 |
| Big_Nose | 74 | 70 | 75 | 74 | 75 | 78 | 81 | **84.56** | 73 | 71 | 79 | 73 | 76 | 81 | **82.02** |
| Bushy_Eyebrow | 80 | 76 | 86 | 80 | 85 | 90 | 88 | **92.80** | 67 | 63 | 79 | 69 | 79 | **82** | 75.24 |
| Chubby | 86 | 82 | 86 | 86 | 86 | 91 | 95 | **95.76** | 67 | 65 | 69 | 68 | 70 | 73 | **74.73** |
| Double_Chin | 88 | 85 | 88 | 90 | 88 | 92 | 96 | **96.42** | 70 | 64 | 75 | 70 | 74 | 78 | **79.17** |
| Narrow_Eyes | 82 | 79 | 84 | 83 | 77 | 81 | 86 | **87.48** | 73 | 68 | 73 | 74 | 77 | **81** | 78.05 |
| Recede_Hair | 76 | 82 | 85 | 84 | 85 | 89 | 92 | **93.82** | 63 | 61 | 84 | 70 | 81 | **85** | 84.11 |
| Earring | 73 | 72 | 78 | 77 | 78 | 82 | 81 | **91.04** | 88 | 85 | 92 | 87 | 90 | **94** | 93.17 |
| Necklace | 68 | 67 | 67 | 70 | 68 | 71 | 87 | **88.15** | 81 | 79 | 86 | 81 | 83 | 88 | **88.81** |
| Average | 81 | 79 | 85 | 83 | 83 | 87 | 87 | **91.66** | 74 | 71 | 81 | 76 | 79 | **84** | 84.34 |

We use 162,770 images as the labelled source data, and 19,962 images as the unlabelled target data on the CelebA dataset. Then, we use FMTNet, which transfers the source data to the target data for facial attribute classification. In this experiment, we randomly select eight labelled attributes ('Attractive', 'Heavy_Makeup', 'No_Beard', 'Oval_Face', 'Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') as the source data, and unlabelled attribute 'Pointy_Nose' as the target data (other attributes can also be used to obtain similar results). Therefore, we build the following eight transfer learning tasks: 'Attractive' → 'Pointy_Nose'; 'Heavy_Makeup'

→ 'Pointy_Nose'; 'No_Beard' → 'Pointy_Nose'; 'Oval_Face' → 'Pointy_Nose'; 'Arched_Eyebrow' → 'Pointy_Nose'; 'Gray_Hair' → 'Pointy_Nose'; 'Male' → 'Pointy_Nose'; 'Smiling' → 'Pointy_Nose'.

As mentioned in Section 3.3, the total loss of the overall TNet is the combination of the softmax loss and three MK-MMD losses. In this experiment, to explore the influence of the correlation between the source and target domains in FMTNet, the loss weight $\alpha$ for the softmax loss (see Eq. (9)) is set to 1.0.

The five attributes ('Attractive', 'Heavy_Makeup', 'No_Beard', 'Oval_Face' and 'Pointy_Nose') are clustered
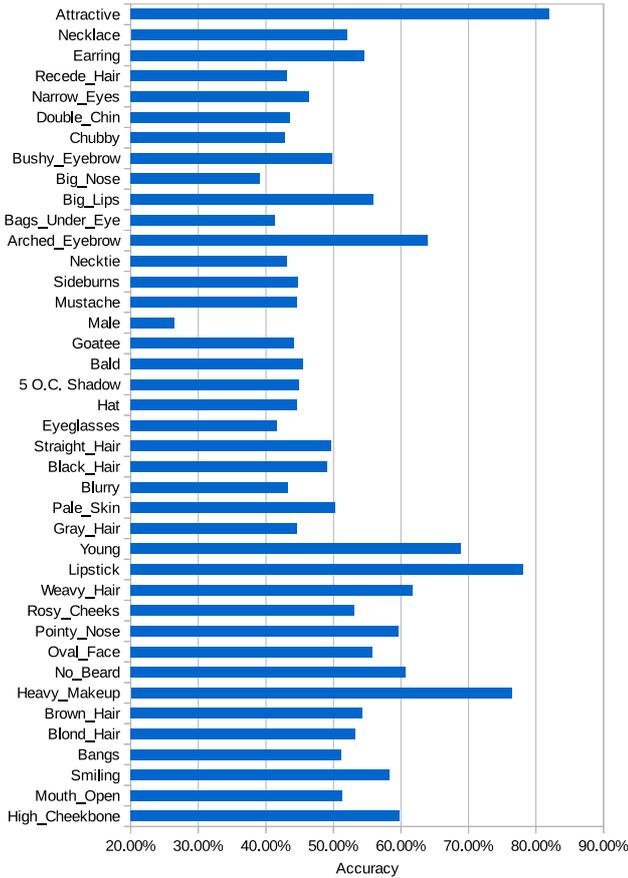
Figure 5: Classification accuracy obtained on different attributes when we directly transfer one attribute to the other attributes. We train one attribute ('Attractive') and directly use the trained model to predict all the 40 attributes on the CelebA dataset.
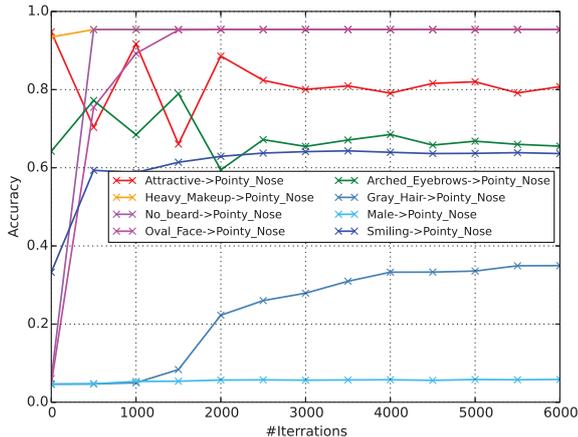


Figure 6: Performance comparison on different source domains transferring to the same target domain during the training process. The source domains have different correlations with the target domain. Note that the loss weight $\alpha$ of the softmax loss in TNet is set to 1.0 in this experiment.

into the same group (see Fig. 1). Therefore, these attributes have high correlation. On the other hand, the four attributes ('Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') and

the 'Pointy_Nose' attribute belong to different groups [11]. In other words, these attributes have low correlation to each other. The performance comparison on different transfer learning tasks over the training process is reported in Fig. 6.

From Fig. 6, the proposed FMTNet obtains better performance for the four transfer tasks ('Attractive' → 'Pointy_Nose', 'Heavy_Makeup' → 'Pointy_Nose', 'No_Beard' → 'Pointy_Nose' and 'Oval_Face' → 'Pointy_Nose') than for the other four transfer tasks ('Arched_Eyebrow' → 'Pointy_Nose', 'Gray_Hair' → 'Pointy_Nose', 'Male' → 'Pointy_Nose' and 'Smiling' → 'Pointy_Nose'). This is mainly because the features from the source attribute and the target attribute are similar when they belong to the same group. Therefore, FMTNet achieves better performance when the correlation between the source attribute and the target attribute is high. In contrast, when the source and target attributes belong to different groups, the performance of the transfer learning task drops because of the significant difference of the data distributions between the two attributes. In summary, the more correlated the two attributes are, the better the transferring performance is (when the value of $\alpha$ is fixed to 1.0).

It is worth pointing that the red curve (i.e., 'Attractive' → 'Pointy_Nose') is not rather stable at the initial 3000 iterations compared with the other curves. This is mainly because that the imbalanced training data are used at the beginning of the whole training process. However, the training accuracy becomes stable as the number of iterations increases, since the classification capability of CNN is gradually enhanced with the number of iterations. Recall that the number of training data for each transfer task is the same. But some facial attributes have many labelled data, while some other facial attributes only have a small number of labelled data. For example, there are totally 4,000 training data for each facial attribute, where 2,500 images are labelled as 'Heavy_Makeup', while only 1,000 images are labelled as 'Attractive'. The small number of labelled data in the source domain (for example, in the task of 'Attractive' → 'Pointy_Nose') will cause that the training accuracy is not stable at the beginning of the training process. In contrast, the training accuracy for some transfer tasks (for example, in the task of 'Heavy_Makeup' → 'Pointy_Nose') is more stable since the labelled data in the source domain are sufficient. Actually, Rudd *et al.* [14] have mentioned the problem of imbalanced data for facial attribute classification.

### 4.3.3. Comparison with the State-of-the-art Methods

In this subsection, we compare the performance of the proposed FMTNet with several state-of-the-art transfer learning methods, including TCA [41], MIDA [42], ITL [43] and GFK [44]. TCA [41] uses the MMD-regularized PCA, which is a conventional transfer learning method, to learn some transfer components across domains. MIDA [42] reduces the discrepancy between the source and target domains by maximizing the independence between the source and target features. ITL [43] consistently learns a domain-invariant feature space and optimizes an information-theoretic metric on the target domain. GFK [44] bridges the source and target domains by interpolating the intermediate subspace.

We use 5,000 images with the labelled information on the CelebA dataset as the source data, and 2,000 images without the labelled information on the CelebA dataset as the target data (also as the test data). The facial attributes are randomly chosen as the source data and the target data. For the traditional transfer learning methods, we extract the features from the source data and target data with single-label learning and FNet, respectively. Then, we use the logistical regression technique to predict the facial attribute.

Firstly, we choose the eight attributes (i.e., 'Attractive', 'Heavy_Makeup', 'No_Beard', 'Oval_Face', 'Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') as the source data and the 'Pointy_Nose' attribute as the target data. According to the experimental results in Section 4.3.2, the performance of FMTNet is significantly affected by the correlation between the source and target attributes. Specifically, when the source and target attributes belong to the same group, FMTNet can achieve excellent results (with $\alpha = 1$). FMTNet uses the MK-MMD loss to measure the difference between the source domain and the target domain, which makes it more discriminative than the other transfer learning methods. However, when the source and target attributes are from different groups, the mean accuracy obtained by FMTNet becomes worse (with $\alpha = 1$). This is mainly due to the fact that when the attributes belong to different groups, the large value of $\alpha$ makes the trained model easily overfit to the source attribute. Therefore, when the source and target attributes belong to different groups, we also set the value of $\alpha$ to be 0.1. In other words, the network focuses more on the minimization of the three MK-MMD losses (mainly reducing the domain discrepancy), when the value of $\alpha$ is small. Therefore, FMTNet uses the MK-MMD loss to measure the difference between the source domain and the target domain, which makes it more discriminative than the other transfer learning methods. The mean accuracy obtained by all the competing methods is shown in Fig. 7, where the loss weight $\alpha$ of the softmax loss is set to 1.0 or 0.1. The 'Same Group' represents the attributes (i.e., 'Attractive', 'Heavy_Makeup', 'No_Beard' and 'Oval_Face') belonging to the same group as the 'Pointy_Nose' attribute. The 'Different Groups' represents the attributes (i.e., 'Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') belonging to different groups compared with the 'Pointy_Nose' attribute.

As shown in Fig. 7, the mean accuracy obtained by FMTNet outperforms that obtained by the competing state-of-the-art methods. Specifically, if the source and target attributes belong to the same group, the mean accuracy obtained by FMTNet achieves the top performance when the value of $\alpha$ is set to 1.0, as shown in Fig. 7 (a). And if the source and target attributes belong to different groups, the mean accuracy obtained by FMTNet achieves the top performance when the value of $\alpha$ is set to 0.1, as shown in Fig. 7 (b). In summary, the loss weight $\alpha$ should be chosen according to the relationship between two attributes for better performance. More specifically, if the source and target attributes belong to the same group, the loss weight $\alpha$ should be assigned with a large value. Otherwise, $\alpha$ should be assigned with a small value.

To show the feature transferability of the proposed method, Fig. 8 demonstrates the t-SNE embeddings [45] of the samples for the 'Attractive' → 'Pointy_Nose' task with the original features and the transferred features (on the source and target domains), respectively. We observe that, with the original features, the data distributions on the source and target domains are quite different. However, with the transferred features obtained by the FMTNet method, the data distributions on the two domains become more similar. This observation further verifies the effectiveness of MK-MMD, which significantly reduces the domain discrepancy.

To further demonstrate the effectiveness of the proposed FMTNet method, we randomly choose 20 facial attributes as the source data and 20 facial attributes as the target data under two different cases (i.e., the attributes in the source domain and in the target domain are from the same group or from different groups). The performance comparison among all the competing methods is given in Table 3. As shown in Table 3, the proposed FMTNet (we set $\alpha = 1$ when the source and target attributes belong to the same group and $\alpha = 0.1$ when the source and target attributes belong to different groups) obtains the top mean accuracy among all the competing methods. Therefore, the loss weight $\alpha$ should be chosen according to the results of attribute grouping (i.e., whether the source and target attributes belong to the same group or not).

The proposed FMTNet method takes advantage of some existing components. However, our method is not simply stacked by these components. In FMTNet, three carefully designed networks (i.e., FNet, MNet and TNet) are used, where these networks share the same structure at their former layers and they differ at their latter layers. Therefore, the networks can be effectively trained via fine-tuning. In other words, TNet can be effectively fine-tuned based on MNet (which is initialized by FNet). Especially, a loss weight scheme is proposed to explicitly exploit the correlation between facial attributes based on attribute grouping, which can improve the generalization performance of the proposed method.

In general, compared with the other transfer learning methods, the proposed FMTNet, which transfers features from the source domain to the target domain, is very effective. The reasons are summarized as follows:

(1) Instead of directly performing transfer learning for facial attribute classification, we take advantage of multi-label learning to learn the network parameters (MNet) and then fine-tune the transfer network (TNet) based on MNet. In other words, we use a hierarchical training strategy, where TNet is fine-tuned based on MNet, and MNet is fine-tuned based on FNet. Therefore, TNet and MNet can effectively trained via fine-tune.

(2) We improve existing multi-label learning and transfer learning components by exploiting the correlation between facial attributes and propose a loss weight scheme based on attribute grouping.

## 5. Conclusions

For the task of facial attribute classification, we have presented a novel method, termed FMTNet, which consists of three
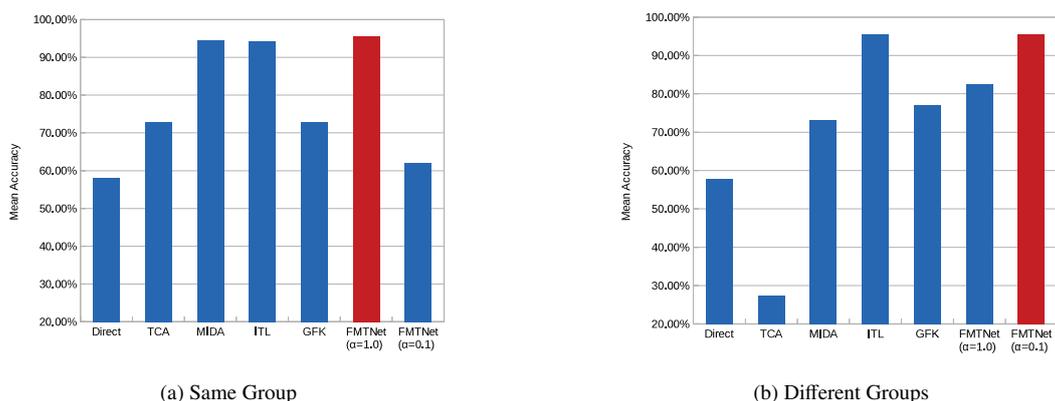
(a) Same Group



(b) Different Groups

Figure 7: The mean accuracy (%) obtained by the proposed FMTNet and five other state-of-the-art methods on the eight attributes (i.e., 'Attractive', 'Heavy_Makeup', 'No_Beard', 'Oval_Face', 'Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') transferring to the 'Pointy_Nose' attribute obtained by the proposed FMTNet and six other state-of-the-art methods. (a) represents the results obtained for the four attributes (i.e., 'Attractive', 'Heavy_Makeup', 'No_Beard' and 'Oval_Face') belonging to the same group as the 'Pointy_Nose' attribute. (b) represents the results obtained for the other four attributes (i.e., 'Arched_Eyebrow', 'Gray_Hair', 'Male' and 'Smiling') belonging to different groups from the one of the 'Pointy_Nose' attribute.



(a) The original features
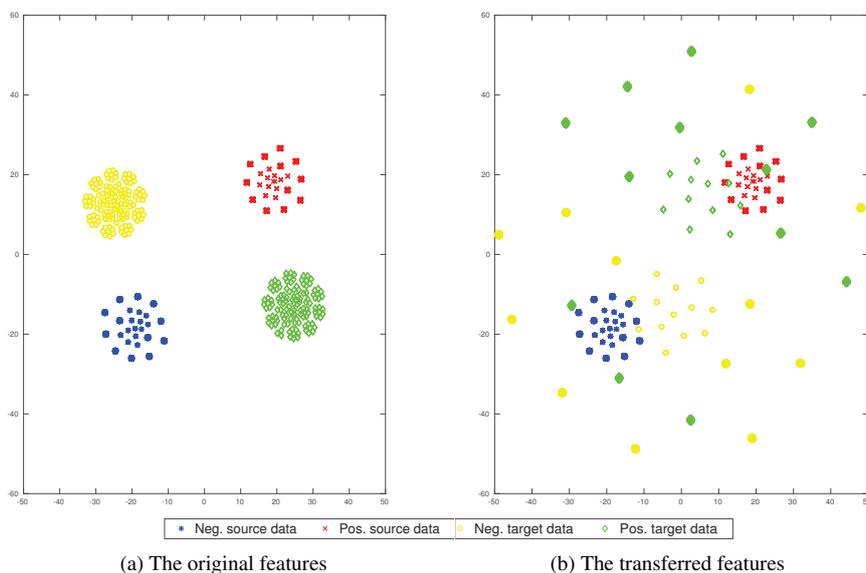


(b) The transferred features

Figure 8: Feature visualization: t-SNE of (a) the original features and (b) the transferred features obtained by FMTNet for 'Attractive' → 'Pointy_Nose' on the source and target domains.

different sub-networks – FNet, MNet and TNet – for face detection, multi-label learning and transfer learning, respectively. MNet in FMTNet predicts multiple facial attributes simultaneously for the labelled facial attributes. Moreover, MNet reduces feature redundancy with the proposed loss weight scheme, resulting in significant performance improvements. Based on MNet, TNet in FMTNet predicts facial attributes with unlabelled information by using MK-MMD. The proposed method, which combines multi-label learning with transfer learning, is general and can be applied to other computer vision tasks.

## References

[1] Y. Sun, X. Wang, X. Tang, Deeply learned face representations are sparse, selective, and robust, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 2892–2900.

14

Table 3: The mean accuracy (%) obtained by the proposed FMTNet in comparison with the state-of-the-art methods. The best results are in boldface.

| Group | Methods | Direct | TCA [41] | MIDA [42] | ITL [43] | GFK [44] | FMTNet |
|---|---|---|---|---|---|---|---|
| Same (α=1.0) | High_Cheekbone→Mouth_Open | 49.65 | 49.75 | 49.85 | 49.95 | 49.95 | **71.95** |
| | Smiling→High_Cheekbone | 35.60 | 49.20 | 49.05 | 49.35 | 49.65 | **74.75** |
| | Attractive→Bangs | 56.65 | 16.00 | 15.85 | 16.85 | 47.75 | **82.90** |
| | Blond_Hair→Brown_Hair | 65.00 | 19.65 | 30.80 | 67.25 | 73.50 | **81.30** |
| | Heavy_Makeup→No_Beard | 72.36 | 82.60 | **85.10** | 81.95 | 28.45 | 75.49 |
| | Pointy_Nose→Oval_Face | 65.20 | 29.45 | 57.75 | 66.65 | 64.25 | **71.15** |
| | Rosy_Cheeks→Wavy_Hair | 37.75 | 38.55 | 61.95 | **62.10** | 59.80 | 62.10 |
| | Lipstick→Young | 75.10 | 75.50 | **76.00** | 75.30 | 40.00 | 75.88 |
| | Gray_Hair→Pale_Skin | 95.90 | 23.45 | 95.95 | 4.45 | 92.75 | **99.05** |
| | Blurry→Gray_Hair | 81.35 | 4.55 | 96.55 | 96.60 | 93.30 | **99.65** |
| | Black_Hair→Straight_Hair | 78.75 | 21.00 | 21.05 | 70.20 | 68.00 | **79.30** |
| | Eyeglasses→ Hat | 5.08 | 12.20 | **95.70** | 95.60 | 93.80 | 95.61 |
| | 5 O.C. Shadow→Bald | 57.51 | 3.45 | 94.95 | **97.60** | 90.45 | **97.60** |
| | Goatee→Male | 36.05 | 39.50 | 61.55 | 60.45 | **63.50** | 61.45 |
| | Mustache→Sideburns | 95.41 | 5.90 | **95.45** | 94.15 | 93.25 | **95.45** |
| | Arched_Eyebrow→Bags_Under_Eye | 30.00 | 20.45 | 40.40 | 65.10 | 73.45 | **80.40** |
| | Big_Lips→Big_Nose | 77.35 | 21.55 | 73.00 | 78.50 | 77.15 | **78.90** |
| | Bushy_Eyebrow→Chubby | 88.18 | 7.20 | 79.65 | 80.70 | 88.20 | **94.15** |
| | Double_Chin→Narrow_Eyes | 21.05 | 27.60 | 85.10 | 85.10 | 82.60 | **85.20** |
| | Recede_Hair→Earring | 21.15 | 21.80 | 78.60 | **79.00** | 75.15 | 79.00 |
| | Average | 57.25 | 28.47 | 67.22 | 68.84 | 70.25 | **82.06** |
| Different (α=0.1) | High_Cheekbone→Attractive | 47.00 | 49.90 | 49.90 | 49.60 | 41.30 | **51.10** |
| | Mouth_Open→Bangs | 61.70 | 15.75 | 15.90 | 16.55 | 44.70 | **84.40** |
| | Smiling→Blond_Hair | 57.85 | 13.10 | 12.55 | 13.25 | 44.35 | **87.05** |
| | Gray_Hair→Brown_Hair | 77.15 | 32.90 | 81.20 | 19.00 | 78.30 | **81.25** |
| | Pale_Skin→Heavy_Makeup | 41.85 | 41.60 | 58.70 | **58.80** | 57.15 | 58.64 |
| | Blurry→No_Beard | 85.90 | 84.70 | 14.00 | 14.15 | 15.45 | **86.04** |
| | Black_Hair→Oval_Face | 67.70 | 29.25 | 29.30 | 64.75 | 60.75 | **71.15** |
| | Straight_Hair→Pointy_Nose | 36.40 | 30.60 | 56.40 | 60.20 | 61.95 | **69.90** |
| | Eyeglasses→Rosy_Cheeks | 8.85 | 15.40 | 92.60 | **92.70** | 89.70 | 90.63 |
| | Hat→Wavy_Hair | 42.95 | 38.55 | 62.05 | 38.55 | 61.05 | **62.10** |
| | 5 O.C. Shadow→Arched_Eyebrow | 30.70 | 30.95 | 67.15 | **69.70** | 63.65 | 68.80 |
| | Bald→Bags_Under_Eye | 78.75 | 26.05 | **80.40** | 20.50 | 79.95 | **80.40** |
| | Goatee→Big_Lips | 54.80 | 35.85 | **65.20** | 64.00 | 64.35 | 64.10 |
| | Male→Big_Nose | 73.90 | 23.50 | 21.70 | 22.50 | 64.05 | **78.86** |
| | Mustache→Bushy_Eyebrow | 60.40 | 14.00 | 87.25 | 86.15 | 85.15 | **87.26** |
| | Sideburns→Chubby | 72.25 | 8.75 | 94.10 | 6.60 | 92.05 | **94.15** |
| | Necktie→Double_Chin | 5.70 | 6.25 | **94.70** | 5.90 | 91.55 | 94.68 |
| | Lipstick→Narrow_Eyes | 48.25 | 15.55 | 15.05 | 15.65 | 56.15 | **63.50** |
| | Young→Recede_Hair | 43.70 | 9.75 | 10.20 | 9.80 | 24.50 | **63.85** |
| | Earring→Bangs | 82.50 | 15.80 | 33.55 | 64.20 | 72.75 | **83.90** |
| | Average | 53.92 | 26.91 | 52.10 | 39.63 | 62.44 | **76.09** |

[2] N. Kumar, A. C. Berg, P. N. Belhumeur, S. K. Nayar, Attribute and simile classifiers for face verification, in: Proceedings of the IEEE International Conference on Computer Vision, 2009, pp. 365–372.

[3] N. Kumar, A. C. Berg, P. N. Belhumeur, S. K. Nayar, Describable visual attributes for face verification and image search, IEEE Trans. Pattern Anal. Mach. Intell. 33 (10) (2011) 1962–1977.

[4] B. Chen, C. Chen, W. H. Hsu, Cross-age reference coding for age-invariant face recognition and retrieval, in: Proceedings of the European Conference on Computer Vision, 2014, pp. 768–783.

[5] Y. F. Yu, D.-Q. Dai, C.-X. Ren, K.-K. Huang, Discriminative multi-layer illumination-robust feature extraction for face recognition, Pattern Recognit. 67 (2017) 201–212.

[6] H. Li, C. Y. Suen, Robust face recognition based on dynamic rank representation, Pattern Recognit. 60 (2016) 13–24.

[7] B. Chen, C. Chen, W. H. Hsu, Face recognition and retrieval using cross-age reference coding with cross-age celebrity dataset, IEEE Trans. Multimed. 17 (6) (2015) 804–815.

[8] N. Cherniavsky, I. Laptev, J. Sivic, A. Zisserman, Semi-supervised learning of facial attributes in video, in: Proceedings of the European Conference on Computer Vision, 2010, pp. 43–56.

[9] P. Luo, X. Wang, X. Tang, A deep sum-product architecture for robust facial attributes analysis, in: Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 2864–2871.

[10] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, L. Bourdev, Panda: Pose aligned networks for deep attribute modeling, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 1637–1644.

[11] Z. Liu, P. Luo, X. Wang, X. Tang, Deep learning face attributes in the wild, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 3730–3738.

[12] S. Kang, D. Lee, C. D. Yoo, Face attribute classification using attribute-aware correlation map and gated convolutional neural networks, in: Proceedings of the IEEE International Conference on Image Processing, 2015, pp. 4922–4926.

[13] M. Ehrlich, T. J. Shields, T. Almaev, M. R. Amer, Facial attributes classification using multi-task representation learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 47–55.

[14] E. M. Rudd, M. Günther, T. E. Boult, Moon: A mixed objective optimization network for the recognition of facial attributes, in: Proceedings of the European Conference on Computer Vision, 2016, pp. 19–35.

[15] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A

large-scale hierarchical image database, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 248–255.

[16] Y. Zhong, J. Sullivan, H. Li, Leveraging mid-level deep representations for predicting face attributes in the wild, in: Proceedings of the IEEE International Conference on Image Processing, 2016, pp. 3239–3243.

[17] P. Luo, Z. Zhu, Z. Liu, X. Wang, X. Tang, Face model compression by distilling knowledge from neurons, in: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, 2016, pp. 3560–3566.

[18] D. G. Lowe, Distinctive image features from scale-invariant key-points, Int. J. Comput. Vis. 60 (2) (2004) 91–110.

[19] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2005, pp. 886–893.

[20] M. L. Zhang, Z. H. Zhou, A review on multi-label learning algorithms, IEEE Trans. Knowl. Data Eng. 26 (8) (2014) 1819–1837.

[21] Y. Huang, W. Wang, L. Wang, T. Tan, Multi-task deep neural network for multi-label learning, in: Proceedings of the IEEE International Conference on Image Processing, 2013, pp. 2897–2900.

[22] I. Pillai, G. Fumera, F. Roli, Designing multi-label classifiers that maximize f measures: State of the art, Pattern Recognit. 61 (2017) 394–404.

[23] B. Wang, J. Tsotsos, Dynamic label propagation for semi-supervised multi-class multi-label classification, Pattern Recognit. 52 (2016) 75–84.

[24] Y. Li, B. Wu, B. Ghanem, Y. Zhao, H. Yao, Q. Ji, Facial action unit recognition under incomplete data based on multi-label learning with missing labels, Pattern Recognit. 60 (2016) 890–900.

[25] S. J. Pan, Q. Yang, A survey on transfer learning, IEEE Trans. Knowl. Data Eng. 22 (10) (2010) 1345–1359.

[26] W. Wei, C. Tian, S. J. Maybank, Y. Zhang, Facial expression transfer method based on frequency analysis, Pattern Recognit. 49 (2016) 115–128.

[27] J. F. Hu, W.-S. Zheng, X. Xie, J. Lai, Sparse transfer for facial shape-from-shading, Pattern Recognit. 68 (2017) 272–285.

[28] M. Long, J. Wang, G. Ding, J. Sun, P. S. Yu, Transfer feature learning with joint distribution adaptation, in: Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 2200–2207.

[29] M. Long, J. Wang, G. Ding, J. Sun, P. S. Yu, Transfer joint matching for unsupervised domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 1410–1417.

[30] J. Yosinski, J. Clune, Y. Bengio, H. Lipson, How transferable are features in deep neural networks?, in: Proceedings of the Neural Information Processing Systems, 2014, pp. 3320–3328.

[31] M. Long, J. Wang, M. I. Jordan, Unsupervised domain adaptation with residual transfer networks, in: Proceedings of the Neural Information Processing Systems, 2016, pp. 136–144.

[32] M. Long, Y. Cao, J. Wang, M. I. Jordan, Learning transferable features with deep adaptation networks, in: Proceedings of the International Conference on Machine Learning, 2015, pp. 97–105.

[33] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv:1409.1556, 2014.

[34] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: Proceedings of the Neural Information Processing Systems, 2015, pp. 91–99.

[35] A. Gretton, B. K. Sriperumbudur, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, Optimal kernel choice for large-scale two-sample tests, in: Proceedings of the Neural Information Processing Systems, 2012, pp. 1205–1213.

[36] Y. Sun, Y. Chen, X. Wang, X. Tang, Deep learning face representation by joint identification-verification, in: Proceedings of the Neural Information Processing Systems, 2014, pp. 1988–1996.

[37] G. B. Huang, M. Ramesh, T. Berg, E. Learned-Miller, Labeled faces in the wild: A database for studying face recognition in unconstrained environments, Tech. rep., Technical Report 07-49, University of Massachusetts, Amherst (2007).

[38] H. Jiang, E. Learned-Miller, Face detection with the faster r-cnn, arXiv:1606.03473, 2016.

[39] N. Kumar, P. Belhumeur, S. Nayar, Facetracer: A search engine for large collections of images with faces, in: Proceedings of the European Conference on Computer Vision, 2008, pp. 340–353.

[40] J. Li, Y. Zhang, Learning surf cascade for fast and accurate object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 3468–3475.

[41] S. J. Pan, I. W. Tsang, J. T. Kwok, Q. Yang, Domain adaptation via transfer component analysis, IEEE Trans. Neur. Netw. 22 (2) (2011) 199–210.

[42] K. Yan, L. Kou, D. Zhang, Domain adaptation via maximum independence of domain features, arXiv:1603.04535, 2016.

[43] Y. Shi, F. Sha, Information-theoretical learning of discriminative clusters for unsupervised domain adaptation, arXiv:1206.6438, 2017.

[44] B. Gong, Y. Shi, F. Sha, K. Grauman, Geodesic flow kernel for unsupervised domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 2066–2073.

[45] J. Donahue, Jia. Y, Vinyals. O, Hoffman. J, Zhang. N, Tzeng. E, Darrell. T, Decaf: A deep convolutional activation feature for generic visual recognition, in: Proceedings of the International Conference on Machine Learning, 2014, pp. 647–655.

[46] C. Xu, T. Liu, D. Tao, C. Xu, Local rademacher complexity for multi-label learning, in: Proceedings of the IEEE International Conference on Image Processing, 2016, pp. 1495–1570.

[47] C. Xu, D. Tao, C. Xu, Robust extreme multi-label learning, in: Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 1275–1284.

[48] V. Jian, N. Modhe, P. Rai, Scalable generative models for multi-label learning with missing labels, in: Proceedings of the International Conference on Machine Learning, 2017, pp. 1636–1644.

[49] Y. Han, F. Wu, Y. Zhang, X. He, Multi-Label transfer learning With sparse representation, IEEE Trans. Circ. & Syst. for Vide. Tech. 20 (8) (2010) 1110–1121.

[50] S. You, C. Xu, Y. Wang, C. Xu, D. Tao, Privileged multi-label learning, in: Proceedings of the International Joint Conference on Artificial Intelligence, 2017, pp. 3336–3342.