# Hierarchical Feature Embedding for Attribute Recognition

Jie Yang[1,2], Jiarou Fan[1], Yiru Wang[1], Yige Wang[1], Weihao Gan[1*], Lin Liu[2], Wei Wu[1]

[1]SenseTime Group Limited, [2]Tsinghua University

takanashiyj@gmail.com,{fanjiarou,wangyiru,ganweihao,wuwei}@sensetime.com,
yige.wang@tum.de,linliu@tsinghua.edu.cn

## Abstract

*Attribute recognition is a crucial but challenging task due to viewpoint changes, illumination variations and appearance diversities, etc. Most of previous work only consider the attribute-level feature embedding, which might perform poorly in complicated heterogeneous conditions. To address this problem, we propose a hierarchical feature embedding (HFE) framework, which learns a fine-grained feature embedding by combining attribute and ID information. In HFE, we maintain the inter-class and intra-class feature embedding simultaneously. Not only samples with the same attribute but also samples with the same ID are gathered more closely, which could restrict the feature embedding of visually hard samples with regard to attributes and improve the robustness to variant conditions. We establish this hierarchical structure by utilizing HFE loss consisted of attribute-level and ID-level constraints. We also introduce an absolute boundary regularization and a dynamic loss weight as supplementary components to help build up the feature embedding. Experiments show that our method achieves the state-of-the-art results on two pedestrian attribute datasets and a facial attribute dataset.*

## 1. Introduction

Attributes, such as gender, hair length, clothing style, are discriminative semantic descriptors that can be used as soft-biometrics in visual surveillance. Attribute recognition concentrates on discerning these attributes of the target human in a given image. It includes pedestrian attribute recognition (PAR), face attribute recognition (FAR), etc. Recently, attribute recognition has received extraordinary attention owing to the potential applications in person re-identification (Re-ID) [22, 26, 36], face verification [21, 6, 46, 34, 45], and human identification [15]. Being a classification problem by nature, it still faces great challenges in real-world scenarios for these reasons: (1) Images
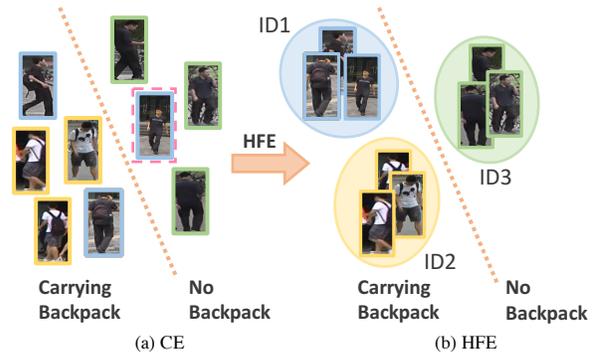


Figure 1. Hierarchical feature embedding on 'backpack' attribute. Images with the same IDs are represented with the same color borders. (a) represents the Cross Entropy feature space, where most of samples can be classified correctly while the hard sample with pink dotted border (whose backpack is totally occluded by the body) is misclassified. (b) With the help of ID constraint, features with the same IDs form fine-grained clusters can pull the hard example back.

might be low-resolution or blurry because of the shot distance or the movements of pedestrians. (2) Different scenes, time slots, angles and poses lead to illumination alterations, viewpoint changes, and appearance variations. (3) Some parts of an object might be occluded by others, resulting in invisibility or ambiguity.

Recently, some methods are introduced to solve these problems and achieve admirable performance. A-AOG [35] explicitly represents the decomposition and articulation of body parts, and accounts for the correlations between poses and attributes. LGNet [28] assigns attribute-specific weights to local features based on the affinity between pre-extracted proposals and attribute locations. These methods aim at applying pivotal parts of images or an attention module to reduce the impact of irrelevant factors to some extent, still they do not cope with the visual appearance variations of attributes as well as occlusion directly. Only attribute-level optimization is focused on in these methods, however the information from attribute recognition related fields, such as person Re-ID, could assist to alleviate vari-

---

*Corresponding author

ation and occlusion issues through imposing stronger constraints.

From the perspective of data, current attribute datasets are all labelled on ID or tracking granularity [25, 32] to reduce workloads. Thus we assume that images captured from the same identity should have the same attributes, but not vise versa. For each attribute, labels are usually coarse-grained owing to the expensive annotation cost. Different persons may get the same attribute labels but with subtle differences in appearance. For example, backpacks with different colors and textures are all tagged as 'backpack'. Therefore, fine-grained feature embeddings are needed for attributes to represent the variety within a class. With identity information, we could set up a two-level feature embedding, i.e., inter-class and intra-class (Fig. 1). For each attribute, samples with identical attribute form coarse-grained classes, while in each attribute class, samples from the same person (with the same attribute definitely) construct the fine-grained ID clusters. We introduce this hierarchical feature embedding model by the following motivations: (1) ID clusters restrict the images with the same ID but variant viewpoints, illumination and appearance to gather more closely, which embed the attribute features with scene invariance and improve robustness. (2) Hard cases for attributes may be easily handled and pulled back by other easy samples of the same ID by the ID constraint, which is difficult to learn only in the attribute level. (3) Like the attribute tags, the ID tags are also utilized in the attribute semantics by holding the assumption above, avoiding integrating different semantic features directly in the same feature space as previous work [25].

Motivated by the above observation, we propose a hierarchical feature embedding (HFE) framework, maintaining the inter-class as well as the intra-class feature embedding by combining attribute and ID information. HFE loss is introduced for fine-grained feature embedding, which consists of two triplet losses and an absolute boundary regularization with the selected quintuplets. With HFE loss constraints, each class could gather more compactly, leading to a more distinct boundary between classes. We propose the absolute boundary regularization for additional absolute constraints because the triplet loss only considers the difference between two distances but ignores the absolute values, and the intra-class triplet loss may interact indirectly with the inter-class boundary. Besides, the quintuplet selection is relevant to the current feature space. However, in the early stage of training, feature spaces are not confident enough for the quintuplet selection, so we design a dynamic loss weight, making the HFE loss weight increasing gradually along with the learning process. In summary, the contributions of this paper are:

- We propose HFE framework to integrate ID information on attribute semantics for fine-grained feature embedding. A novel HFE loss is introduced for both inter-class and intra-class level constraints.

- We construct an absolute boundary regularization by strengthening the original triplet loss with an absolute constraint.

- We introduce a dynamic loss weight, which forces the feature space to transit from origin to the improved HFE-restricted space by degrees.

- The proposed method is evaluated on two pedestrian attribute datasets and one face attribute dataset. Experiments show our method achieves the state-of-the-art results on all three datasets.

## 2. Related Work

### 2.1. Attribute Recognition

Recently, deep learning based attribute recognition methods achieve impressive performance. In PAR, these methods includes global based [23, 41, 1, 7], local parts based [28], visual attention based [31], sequential prediction based [55] methods, etc. Among them, DeepMar [23] is an early global based PAR work. Considering the imbalanced data distribution, it proposes cost-sensitive cross entropy (CE) loss for classification. It also proposes a new loss to handle imbalanced data as well as a new attention mechanism. LGNet [28] assigns attribute-specific weights to local features based on the affinity between pre-extracted proposals and attribute locations. Hydraplus-Net [31] proposes an attention based model and exploits the global and local contents with multi-level feature fusion of a single pedestrian image. ALM [47] aims at learning attribute-aware representation through attribute localization. Attribute-aware attention model [11] exploits the correlation between global features and attribute-specific features and utilize it to generate attention mask in a reciprocal manner. Localizing by describing [30] learns the bounding boxes related to the location of attributes explicitly using REINFORCE algorithm with a designed reward function in a weakly supervised manner. [39] designs an attention mechanism for aggregating multi-scale features as well as a loss function similar to focal loss [24] in order to tackle the imbalanced data problem. GRL [55] proposes a RNN based grouping recurrent learning method that exploits the intra-group mutual exclusion and inter-group correlation. FAR can be also divided into part-based and holistic approaches [38, 12].

ReID aims at matching a target person in a set of query pedestrian images. A great amount of deep learning based ReID works provide promising solutions [48, 8, 4, 2, 9, 37]. A bunch number of existing methods rely on exploiting discriminative features, which is in the same spirit as fine-grained recognition. Attributes and ReID are highly correlated pedestrian visual appearance representations, but vary
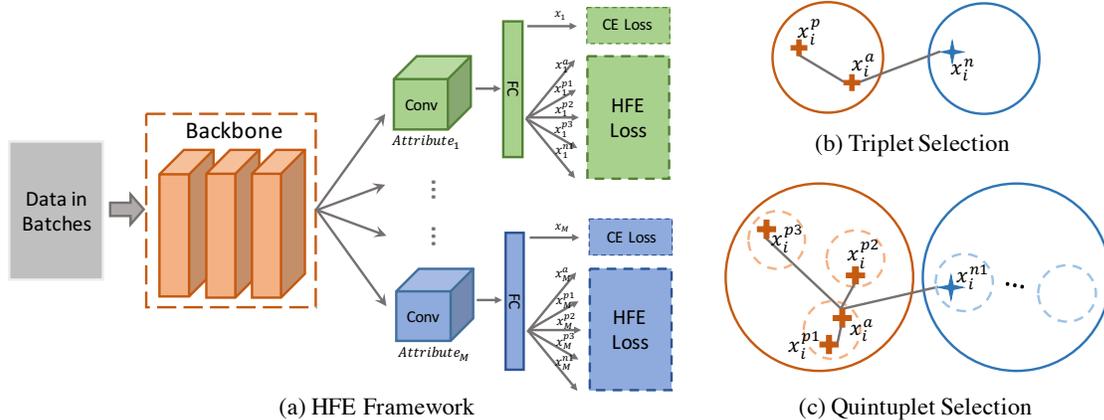
Figure 2. (a) Overview of the proposed hierarchical feature embedding (HFE) framework, it consists of a backbone model, attached by $M$ branches for $M$ attributes. We calculate the CE loss and HFE loss based on the quintuplet selection in each branch. (b) and (c) are triplet and quintuplet selection respectively, where orange and blue represent the different attribute classes.

in semantics and granularities. Even though they are different tasks, the common characteristic can be beneficial to both of them, which is exploiting discriminative features. And therefore they are reasonable to be dealt with jointly. Some works utilize both information for multi-task learning [42] or assisting the main task [25, 27]. The methods can be summarized in two categories: (1) shared backbone and task-independent branches (2) task-independent models and combining high level features in some way (e.g., concatenated FC). For example, APR [25] learns ReID embedding and pedestrian attributes simultaneously, sharing the same backbone and owning classification FC layers respectively. UF [42] trains two different models for two tasks and concatenates branches to one identity vector for ReID. These methods combine these two kind of features to some extent. However incorporating them into the coefficient feature representation indiscriminately is less powerful since attribute recognition and ReID are essentially divergent. Persons with similar attributes can also be different identities. Therefore a more rational way of combining both information is required.

## 2.2. Metric Learning

The objective of metric learning [54, 17] is to learn a proper metric feature space so that the distances between similar samples reduce and that of dissimilar samples enlarge. While traditional metric learning algorithms [19] are based on linear transformation, nonlinear models, etc, due to the recent advances in deep learning, Convolutional Neural Networks have been a powerful tool to learn a task-specific metric and achieved impressive results in a wide range of tasks. Many metric learning algorithms have been proposed in image retrieval [50], ReID [5], face recognition [40, 44, 29, 53], etc. Representative methods are contrastive loss and triplet loss. Contrastive loss [10] restricts

pair inputs and results in distances between similar pairs as close as possible and that of dissimilar pairs to be larger than margin. Triplet loss [40] applies triplet as input and ensures the difference between the distance of (anchor, negative) feature and (anchor, positive) feature is larger than margin. Beyond triplet loss, quadruplet loss [5] and quintuplet loss [13] are also introduced to improve performance. Center loss [52] is designed for face recognition strives to push samples to their respective clusters centers. We propose HFE loss by applying inter-class and intra-class triplet losses for fine-grained constraints.

## 3. Proposed Method

**Problem Definition.** Given $N$ images $\{I_1, I_2, ..., I_N\}$ and each image $I_j$ has $M$ visual attribute tags $y_j = \{y_{j1}, y_{j2}, ..., y_{jM}\}$ together with ReID label $l_j$. The images from the same person are labelled with identical attributes, i.e., $\forall_{l_i=l_j} y_i = y_j$. ReID auxiliary attribute recognition aims at training a model containing the attribute and ID information to predict the attributes $y_k$ for the characteristic of the person in an unseen image $I_k$.

**Network Architecture.** As shown in Fig. 2 (a), the proposed hierarchical feature embedding (HFE) network consists of a backbone model, to which $M$ branches for $M$ attributes are attached. In the shared backbone, the model learns a common feature embedding for all attributes. For each attribute, we construct branches respectively for two reasons: (1) Different attributes, such as age and gender, should own their specific feature embeddings. (2) We construct metric loss for each attribute in their own feature spaces, which can not be applied on a shared feature space. For example, there are two images $I_1, I_2$ from different IDs, and the attributes are (long hair, carrying backpack) and (long hair, no backpack). We should pull them closer for hair length feature while push them away for backpack

feature. Each attribute branch contains Conv-BN-ReLU-Pooling-FC sequential layers. We calculate the Cross Entropy (CE) loss and metric loss on each branch.

**Loss Computation.** We apply CE loss for attribute classification (Eq. 1) as most works do. Besides, an HFE loss is utilized for auxiliary metric learning with weight $w$ (Eq. 2). HFE loss consists of inter-triplet loss, intra-triplet loss and Absolute Boundary Regularization, which will be introduced in the next section.

$$L_{CE} = -\frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M} y_{ij}log(p_{ij}) + (1-y_{ij})log(1-p_{ij}) \tag{1}$$

$$Loss = L_{CE} + wL_{HFE} \tag{2}$$

### 3.1. Hierarchical Feature Embedding

**Triplet Loss**. Triplet loss has been widely used for metric learning. As shown in Fig. 2 (b), it trains on a series of triplets $\{x_i^a, x_i^p, x_i^n\}$, where $x_i^a$ and $x_i^p$ are image features from the same label, and $x_i^n$ from a different label. $a$, $p$ and $n$ are abbreviations of anchor, positive and negative sample respectively. The formulation is as follows:

$$L_{trp} = \frac{1}{N}\sum_{i=1}^{N}[d(x_i^a, x_i^p) - d(x_i^a, x_i^n) + \alpha]_+ \tag{3}$$

Here, $d(.)$ represents the Euclidean distance, and $\alpha$ is the margin which forces the gap of $d(x_i^a, x_i^n)$ and $d(x_i^a, x_i^p)$ larger than $\alpha$. $[z]_+$ means $max(z, 0)$. When the gap is larger than $\alpha$, the triplet loss would be zero.

**Inter-class Triplet Loss.** We can extend triplet loss to attribute classification scenario. As shown in Eq. 4, $x_{ij}^a$ is the feature of anchor sample $I_i$ of $j$-th attribute, associated with the feature of a positive and negative sample in regard to $x_{ij}^a$, i.e, $x_{ij}^{p3}, x_{ij}^{n1}$. Here we define the triplet on the attribute level, and $\alpha_1$ is the inter-class margin.

$$L_{inter} = \frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M}[d(x_{ij}^a, x_{ij}^{p3}) - d(x_{ij}^a, x_{ij}^{n_1}) + \alpha_1]_+$$
$$y_{ij}^{p3} = y_{ij}^a, \ y_{ij}^{n_1} \neq y_{ij}^a, \ l_{ij}^{p3} \neq l_{ij}^a, \ l_{ij}^{n_1} \neq l_{ij}^a \tag{4}$$

We use batch hard mode [40] for triplet selection. In each batch, we take the closest negative sample to anchor as the hard negative sample $x_{ij}^{n_1} = argmin_{x_{ij}^n}d(x_{ij}^a, x_{ij}^n)$, and the farthest positive as the hard positive sample $x_{ij}^{p3} = argmax_{x_{ij}^p}d(x_{ij}^a, x_{ij}^p)$, which is called as the hard inter-class triplet loss.

| item | attribute | ID | distance |
|------|-----------|-----|----------|
| $x_{ij}^{p_1}$ | the same | the same | farthest |
| $x_{ij}^{p_2}$ | the same | different | nearest |
| $x_{ij}^{p_3}$ | the same | different | farthest |
| $x_{ij}^{n}$ | different | different | nearest |

Table 1. The summary of quintuplet. $x_{ij}^a$ is the anchor and not listed above.

**Intra-class Triplet Loss.** With the attribute inter-class triplet loss, we could separate the feature embedding between classes. However, the feature embeddings in each class are still mixed up. Intuitively, the features of samples with similar appearance or the same ID should be closer than others. However, it is not easy to get such a perfect feature embedding without extra constraint on the intra-class level. To form ordered and fine-grained intra-class feature embeddings, we utilize ID information to enforce the features that belong to the same person gathered more closely. We construct the intra-class feature embedding for these two reasons: (1) The intra-class triplet loss restricts the features from the same person to gather more closely, making the embedding more robust to scene variance. (2) The hard cases for attributes but not for ID can be easily handled by ID clusters in the attribute feature embedding. Here we introduce a hard intra-class triplet loss, similar to the hard inter-class triplet loss, while the hard negative sample is replaced by the closest positive sample to the anchor with different ID but the same attribute, $x_{ij}^{p_2} = argmin_{x_{ij}^p}d(x_{ij}^a, x_{ij}^p)$ for $y_{ij}^p = y_{ij}^a$, $l_{ij}^p \neq l_{ij}^a$, and the hard positive sample is turned into the farthest positive sample with the same ID (with the same attribute definitely), $x_{ij}^{p_1} = argmax_{x_{ij}^p}d(x_{ij}^a, x_{ij}^p)$ for $y_{ij}^p = y_{ij}^a$, $l_{ij}^p = l_{ij}^a$. The intra-class triplet loss is shown in Eq. 5, and $\alpha_2$ is the intra-class margin.

$$L_{intra} = \frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M}[d(x_{ij}^a, x_{ij}^{p_1}) - d(x_{ij}^a, x_{ij}^{p_2}) + \alpha_2]_+$$
$$y_{ij}^{p_1} = y_{ij}^a, \ y_{ij}^{p_2} = y_{ij}^a, \ l_{ij}^{p_1} = l_{ij}^a, \ l_{ij}^{p_2} \neq l_{ij}^a \tag{5}$$

To maintain the structures of inter-class and intra-class feature embedding simultaneously, we incorporate the inter-class and intra-class triplet loss. As shown in Fig. 2 (c), HFE loss takes quintuplet samples $\{x_{ij}^a, x_{ij}^{p_1}, x_{ij}^{p_2}, x_{ij}^{p_3}, x_{ij}^{n_1}\}$ as input and manage to maintain the multiple relative relationships, $d(x_{ij}^a, x_{ij}^{p_1}) < d(x_{ij}^a, x_{ij}^{p_2}) < d(x_{ij}^a, x_{ij}^{p_3}) < d(x_{ij}^a, x_{ij}^{n})$. The quintuplet is summarized in Table 1. With the constraints on both inter-class and intra-class level, we can construct a hierarchical feature embedding to incorporate attribute information as well as ID information in the attribute feature space.

## 3.2. Absolute Boundary Regularization

Triplet loss only restricts the difference between $d(x_{ij}^a, x_{ij}^p)$ and $d(x_{ij}^a, x_{ij}^n)$ while ignores the absolute value. The difference is dependent on the selected triplet in batch, which is hard to ensure $d(x_{ij}^a, x_{ij}^p) < d(x_{ij}^a, x_{ij}^n)$ in the whole training dataset.

In our constraints, to guarantee attributes owing discriminative intra-class feature embeddings, we pull away $x_{ij}^{p_2}$ from $x_{ij}^a$ relative to $x_{ij}^{p_1}$. Although margin $\alpha_2 < \alpha_1$, $L_{intra}$ may interact indirectly with inter-class boundary.

Considering these two factors, we force $d(x_{ij}^a, x_{ij}^n)$ to be larger than an absolute distance $\alpha_3$, named as Absolute Boundary Regularization (ABR).

$$L_{ABR} = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} [\alpha_3 - d(x_{ij}^a, x_{ij}^n)]_+ \qquad (6)$$

We compose these three parts in Eq. 7 as HFE loss, which aims at a hierarchical feature embedding, considering not only inter-class features but also inner-class, not only relative distance but also absolute to ensure the boundary more discriminative.

$$L_{HFE} = L_{inter} + L_{intra} + L_{ABR} \qquad (7)$$

## 3.3. Dynamic Loss Weight

In the early stage of training, the feature space is not good enough for quintuplet selection. Therefore, applying a large weight for HFE loss at the beginning may amplify the noise caused by the initial model. To solve this problem, we set a small weight for HFE loss at the beginning, making the original CE loss play a major role in optimization and produce an elementary feature space. Then we enlarge the weight afterwards to refine the original feature space to be more fine-grained. So we introduce a dynamic loss weight, modifying the composite function proposed by [51] to control the weight increasing from small to large nonlinearly. In Eq. 8, $T$ means the total training iterations, and $iter$ means the current iteration. $w_0$ is a given value.

$$w = [\frac{1}{2}cos(\frac{T - iter}{T}\pi) + \frac{1}{2}]w_0 \qquad (8)$$

In the training process, we increase the HFE loss weight from zero to the given value, forcing the feature space to transit from origin to the improved HFE-restricted space by degrees.

## 4. Experiment

### 4.1. Datasets

**Market 1501 attribute dataset** [25] is an extension of Market-1501 dataset [56] with person attribute annotations.

It contains 32,668 annotated bounding boxes of 1,501 identities and 12 different types of attribute annotations for each identity. Attributes include 10 binary attribute (such as gender, hair length and sleeve length) and 3 multi-class attributes (i.e., age, upper clothing color and lower clothing color). Images are captured from six cameras and each annotated identity is present in at least two cameras.

**Duke attribute dataset** [25] is an extension of DukeMTMC-ReID dataset [57] with person attribute annotations. It contains 36,411 bounding boxes of 1,404 identities and 10 different types of attribute annotations for each identity. Attributes includes 8 binary attribute (such as gender, length of upper-body clothing and wearing boots) and 2 multi-class color attributes for upper clothing and lower clothing. Images are captured from eight cameras and each annotated identity is present in at least two cameras.

**CelebA** [32] is a large-scale face attribute dataset with annotations of 40 binary classifications (such as eyeglasses, bangs and pointly nose). The dataset contains 202,599 images from 10,177 identities and covers large pose variations and background clutter.

### 4.2. Evaluation

For the first two PAR datasets, we evaluate attribute recognition performance on both class-based and instance-based level. (1) **Class-based:** We calculate the classification accuracy for each attribute class and report the mean accuracy of all attributes [25]. (2) **Instance-based:** We measure the accuracy, precision, recall and F1 score for all test samples. For accuracy, precision and recall, we first compute the scores of predicted attributes against the gound truth for each test sample image and then get the average scores over all test cases. The F1 score is computed based on precision and recall [23]. The gallery images are used as test set and we transform the multi-class attributes into binary classes [25].

For face attribute dataset, we evaluate class-based mean accuracy. [12]

### 4.3. Implementation Details

The common settings for PAR and FAR: We use Adam [18] as an optimization algorithm. The weight decay is 5e-4, batch size is 256. We random sample 64 identities with 4 images for each to form a batch. Horizontal flip is applied during the training process. We set $\alpha1$, $\alpha2$, $\alpha3$ and $w_0$ with 0.3, 0.1, 5, and 1 respectively.

For PAR, we exploit ResNet 50 as the backbone. The base learning rate is 2e-4 and decays exponentially after epoch 50. We train 130 epochs in total.

For FAR, we exploit DeepID2 [43] as the backbone. To accommodate with backbone, in each attribute branch, we replace the convolutional layers with linear layers. The base learning rate is 1e-2 and we use cosine annealing schedule

| Method | gender | hair | L.slv | L.low | S.clth | hat | B.pack | bag | H.bag | age | C.up | C.low | avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CE | 84.86 | 84.51 | 93.49 | 90.64 | 90.24 | 96.01 | 79.58 | 74.16 | 86.67 | 92.62 | 92.99 | 92.86 | 88.22 |
| DeepMAR | 85.24 | 85.48 | 92.79 | 91.37 | 89.37 | 95.93 | 84.56 | 71.57 | 86.53 | **94.56** | 93.78 | 92.92 | 88.68 |
| APR | 85.64 | 85.62 | 92.87 | 92.80 | 89.91 | 97.13 | 78.12 | 75.41 | <u>90.53</u> | 93.18 | 93.18 | 92.77 | 88.93 |
| UF | 88.94 | 78.26 | 93.53 | 92.11 | 84.79 | 97.06 | 85.46 | 67.28 | 88.4 | 84.76 | 87.5 | 87.21 | 86.28 |
| JCM | 89.7 | 82.5 | <u>93.7</u> | **93.3** | 89.2 | <u>97.2</u> | 85.2 | **86.9** | 86.2 | 87.4 | 92.4 | 93.1 | 89.73 |
| HP Net | <u>94.74</u> | 89.11 | 93.14 | 92.47 | 92.06 | 96.94 | <u>87.33</u> | 79.65 | 89.22 | 93.28 | 93.21 | 92.39 | 91.13 |
| DIAC | 93.32 | <u>90.43</u> | 93.24 | 92.82 | **95.63** | 96.98 | 86.18 | 77.56 | 88.84 | 92.88 | **94.55** | <u>92.84</u> | <u>91.27</u> |
| HFE (Ours) | **94.88** | **90.51** | **94.03** | <u>93.25</u> | <u>94.18</u> | **97.88** | **90.41** | <u>85.35</u> | **91.45** | <u>94.37</u> | <u>94.43</u> | **94.00** | **92.90** |

Table 2. Class-based evaluation on Market 1501 attribute dataset with the best results in bold and the second best results underlined. 'L.slv', 'L.low', 'S.clth', 'B.pack', 'H.bag', 'C.up', 'C.low' denote length of sleeve, length of lower-body clothing, style of clothing, backpack, handbag, color of upper-body clothing and color of lower-body clothing, respectively.

| Method | gender | L.up | boots | hat | B.pack | bag | H.bag | C.shoes | C.up | C.low | avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CE | 82.33 | 86.63 | 88.36 | 82.98 | 73.31 | 80.65 | 91.60 | 90.92 | 95.40 | 91.45 | 86.36 |
| DeepMAR | 82.26 | 87.14 | 88.49 | 82.15 | 75.84 | 82.54 | 91.53 | <u>91.38</u> | 95.18 | 92.52 | 86.90 |
| APR | 83.47 | 87.44 | 88.02 | 86.98 | 75.79 | 82.16 | 92.61 | 90.67 | 94.23 | <u>97.43</u> | 87.88 |
| UF | **88.94** | **93.6** | 80.13 | 82.97 | 87.02 | <u>91.60</u> | 89.60 | 83.65 | 93.94 | 91.84 | 88.33 |
| JCM | <u>87.4</u> | 88.3 | 89.6 | 83.3 | **89.0** | 87.9 | 92.4 | 87.1 | 92.9 | 92.1 | 89.00 |
| HP Net | 83.91 | 87.58 | 86.72 | 78.91 | 77.54 | 83.37 | 93.40 | 88.91 | **96.81** | 97.19 | 87.43 |
| DIAC | 85.87 | 89.74 | <u>89.63</u> | **90.79** | 82.90 | 87.88 | <u>93.47</u> | 90.21 | <u>95.92</u> | 97.11 | <u>90.35</u> |
| HFE (Ours) | 87.02 | <u>89.88</u> | **90.70** | <u>88.69</u> | <u>88.50</u> | **91.81** | **93.64** | **93.82** | 95.85 | **97.80** | **91.77** |

Table 3. Class-based evaluation on Duke attribute dataset with the best results in bold and the second best results underlined. 'L.up', 'B.pack', 'H.bag', 'C.shoes', 'C.up', 'C.low' denote length of sleeve, backpack, handbag, color of shoes, color of upper-body clothing and color of lower-body clothing, respectively.

[33]. We train 300 epochs in total.

## 4.4. Experiments on Pedestrian Attribute Dataset

We list the results of state-of-the-art methods on these two pedestrian attribute datasets, i.e., Market 1501 attribute dataset and Duke attribute dataset. Table 2 and Table 3 show class-based metrics and Table 4 indicate instance-based evaluations. Among the compared methods, CE means applying CE loss for attribute classification with the same backbone as ours. DeepMAR[23] applies weighted CE loss function. APR[25], UF [42] and JCM[27] are three methods combining attibute and ID information. HP Net[31] and DIAC [39] are attribute-focused methods that achieve competitive performance recently.

With 8 Titan XP GPUs, our method costs 73.6 seconds per epoch on average for Duke dataset while the CE loss costs 61.2 seconds. Extra time consumption is only on training phase, and no additional computation is needed for inferring.

With regard to attribute and ID joint methods, performances are improved to some extent, but are still inferior to CE for some attributes, such as 'backpack', 'boots' for APR as well as 'hair', 'color of shoes' for UF and JCM. It shows that combining attribute recognition and ReID in a shared feature space directly is harmful for some attributes. These attributes may not contribute to ReID, resulting in fewer attention by combining ReID directly. However, our method utilizes the ID tags in the semantics of attributes to build up a better feature embedding for attributes, therefore surpassing previous joint methods significantly. Compared with attribute-focused methods, our method achieves better results on highly identity-related attributes, such as 'gender' and 'age', as well as variant and subtle attributes, such as 'bag' and 'boots', with the extra ID information.

Overall, Our method achieves the best performance on both datasets in five evaluation metrics, outperforming the second best results by 1.63%, 2.98%, 1.77%, 2.34%, 2.13% on Market 1501 attribute dataset, and 1.42%, 2.66%, 1.52%, 0.96%, 2.23% on Duke attribute dataset in average accuracy, accuracy, precision, recall and F1 respectively.

For a more intuitive analysis, we demonstrate the recognition results for two IDs with three images each in Fig. 3. HFE achieves better performance for most attributes, especially for accessories such as bag and backpack, which are sensitive to angle and pose variations. When the attribute object is clearly visible, all three methods achieve good performance. However, when occlusion happens, HFE still predicts correctly thanks to the ID constraint, while the other two methods perform rather poorly.

## 4.5. Ablation Study

The advantage of HFE is its capability of learning a fine-grained comprehensive attribute feature representation involving identity information. To better illustrate this, we

| Metric | Market 1501 attribute dataset | | | | Duke attribute dataset | | | |
|---|---|---|---|---|---|---|---|---|
| Method | acc. | prec. | recall | F1 | acc. | prec. | recall | F1 |
| CE | 69.21 | 82.55 | 78.43 | 80.44 | 69.00 | 81.39 | 77.24 | 79.26 |
| DeepMAR | 69.65 | 82.60 | 80.24 | 81.40 | 70.67 | 81.82 | 82.24 | 82.03 |
| APR | 70.25 | 83.52 | 78.96 | 81.18 | 70.10 | 82.74 | 79.02 | 80.83 |
| HP Net | 74.82 | 85.26 | 83.31 | 84.27 | 67.63 | 82.77 | 75.19 | 79.79 |
| DIAC | 75.03 | 85.64 | 83.18 | 84.39 | 74.02 | 84.85 | 83.44 | 83.14 |
| HFE | **78.01** | **87.41** | **85.65** | **86.52** | **76.68** | **86.37** | **84.40** | **85.37** |

Table 4. Instance-based evaluation on Market 1501 attribute dataset and Duke attribute dataset with the best results in bold and the second best results underlined.



(a) Bag       (b) Backpack

Figure 3. Visualization of recognition results. Correct predictions are bounded by green box and red box otherwise. (a) and (b) are prediction results of two IDs on bag and backpack respectively, each compared by three methods.

| Metric Loss | avg. | acc. | prec. | recall | F1 |
|---|---|---|---|---|---|
| $None$ | 88.82 | 70.03 | 83.15 | 79.27 | 81.12 |
| $L_{inter}$ | 90.83 | 71.83 | 85.35 | 80.63 | 82.93 |
| $L_{intra}$ | 91.27 | 74.61 | 85.43 | 83.56 | 84.48 |
| $L_{inter} + L_{intra}$ | 92.44 | 77.08 | 86.73 | 84.88 | 85.99 |
| $L_{inter} + L_{intra} + ABR$ | 92.73 | 77.57 | 87.00 | 85.45 | 86.22 |
| $L_{inter} + L_{intra} + ABR*$ | **92.90** | **78.01** | **87.41** | **85.65** | **86.52** |
| $L_{inter} + L_{pairwise\_intra}$ | 92.30 | 76.49 | 86.21 | 84.53 | 85.36 |

Table 5. Ablation study on Market 1501 attribute dataset. * means replacing the fixed loss weight with dynamic setting.

analyze the effectiveness of each component with quantitative comparisons and qualitative visualization. The ablation study is done on market 1501 attribute dataset.

**Quantitative Evaluation.** Table 5 quantifies the benefits of our hierarchical feature embedding, absolute boundary regularization (ABR), and dynamic loss weight respectively. Based on the CE method, the performance improves distinctly with only $L_{inter}$. While for $L_{intra}$, the improvement is more significant than $L_{inter}$, indicating the extra ID clusters assist attribute classification indeed. Com-

bining $L_{inter}$ and $L_{intra}$ achieves better performance than each of them separately, demonstrating the complementarity of $L_{inter}$ and $L_{intra}$. The combination achieves 3.62%, 7.05% and 4.87% compared with the CE loss method for avg., acc. and F1 respectively. The improvements of ABR and the dynamic loss weight are also demonstrated in table 5. Finally, with combining all the components, HFE achieves 4.08%, 7.98% and 5.40% improvements for avg., acc. and F1 respectively in total.

In order to verify the necessity of setting up intra-class discriminative embedding by the triplet loss, we replace $L_{inter}$ with a simple pairwise loss which only considers IDs' intra-class compactness but no separability with different IDs. As a result, the triplet loss achieves sightly better performance, which means the inter-class separability of ID is helpful for a fined-grained feature embedding and maintaining detailed attribute information.

**Qualitative Evaluation.** We proceed with qualitative evaluation in both inter-class and intra-class level. Fig. 4 demonstrates the learnt feature embedding visualization of attribute 'gender' for different loss function. The CE
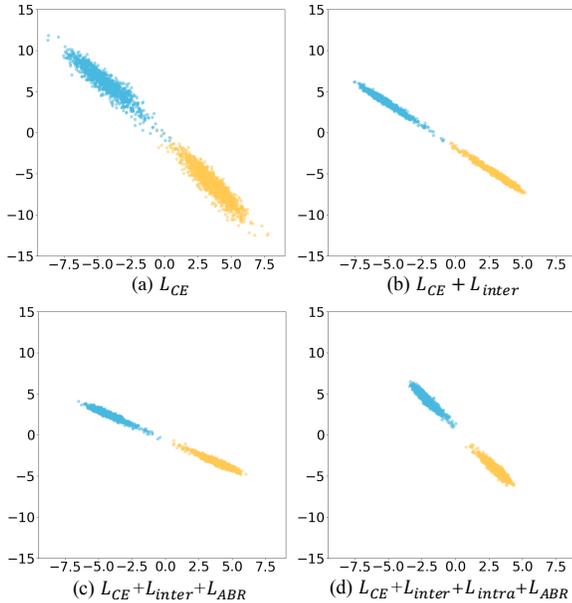
Figure 4. Feature embedding visualization of attribute 'gender'. Blue points represent 'female' and yellow points represent 'male'. Figures are in the same scale for fair comparison.
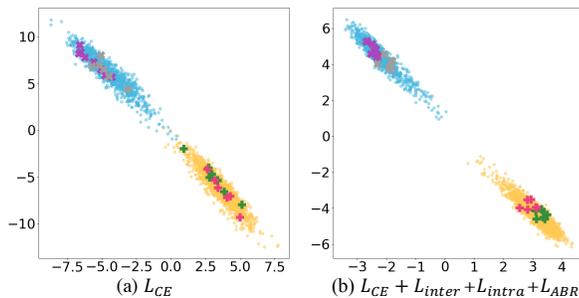


Figure 5. ID clusters are visualized for two methods. Four IDs are shown. '+' and 'x' markers with the same colors come from the identical ID.

loss produces two attribute clusters but the boundary is still in mess. Obviously, applying $L_{inter}$ makes the clusters more compact and the margin between classes more distinct. Besides, adding the constraint of $L_{ABR}$, the effect is more prominent. Furthermore, with adding $L_{intra}$, the two classes are restricted by both inter-class and intra-class constraints, leading to tighter clusters and the most discriminating boundary.

Besides, we evaluate the ID clustering ability and effect. As Fig. 5 shows, CE does not control the intra-class arrangements so the IDs are disorganized, whereas HFE can form fine-grained intra-class clusters by the constraint of $L_{intra}$. We can see that from introducing the intra-class constraint, each class gathers more closely and the margin between them is much more distinct, so hierarchical feature embedding does help to classify attributes.

| Method | acc. |
|---|---|
| FaceTracer | 81.12 |
| PANDA-l | 85.00 |
| LNets+ANet | 87.30 |
| Walk-and-Learn | 88.00 |
| Rudd et al. Moon | 90.94 |
| CLMLE | 91.13 |
| SSP + SSG | 91.80 |
| HSAI | 91.81 |
| HFE | 92.17 |

Table 6. Class-based accuracy on CelebA.

## 4.6. Experiments on Face Attribute Dataset

To evaluate the generalization ability of our framework, we conduct the experiment on a face attribute dataset CelebA. Other state-of-the-art methods are FaceTracer [20], PANDA-l [3], LNets+ANet [32], Walk-and-Learn [49], Rudd et al. Moon [38], CLMLE [14], SSP + SSG [16] and HSAI [12] respectively. Table 6 shows HFE achieves 92.1% accuracy and outperforms all other methods. With the help of face ID information, HFE could also build up a fine-grained feature embedding for face attributes, indicating our HFE framework can be easily generalized to similar scenarios, providing a general framework for fine-grained recognition.

## 5. Conclusion

In this paper, we present a novel end-to-end Hierarchical Feature Embedding (HFE) framework to explore the combination of attribute and ID information in attribute semantics for attribute recognition. In HFE, each attribute-class is discriminative by the inter-class constraint. Moreover, with the supplementary ID information, we maintain the ID clusters in each attribute class by the intra-class constraint for a fine-grained feature embedding. We apply ID information in attribute semantics and refrain from combining attribute and ID information in the same feature space directly. Furthermore, our mechanism introduce a coarse-to-fine process for discriminative fine-grained feature learning. In addition, we introduce an Absolute Boundary Regularization for combining relative and absolute distance constraint. We also design a dynamic loss weight to force the feature space transiting from the origin to the improved HFE-restricted space by degrees, facilitating the performance of our model and the stability of training. Extensive ablation studies and experimental evaluations justify effectiveness of our proposed method.

## 6. Acknowlegements

# References

[1] Abrar H Abdulnabi, Gang Wang, Jiwen Lu, and Kui Jia. Multi-task cnn model for attribute prediction. *IEEE Transactions on Multimedia*, 17(11):1949–1959, 2015. 2

[2] Ejaz Ahmed, Michael Jones, and Tim K Marks. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3908–3916, 2015. 2

[3] Lubomir Dimitrov Bourdev. Pose-aligned networks for deep attribute modeling, July 26 2016. US Patent 9,400,925. 8

[4] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. Abdnet: Attentive but diverse person re-identification. *arXiv preprint arXiv:1908.01114*, 2019. 2

[5] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 403–412, 2017. 3

[6] Sumit Chopra, Raia Hadsell, Yann LeCun, et al. Learning a similarity metric discriminatively, with application to face verification. In *CVPR (1)*, pages 539–546, 2005. 1

[7] Yubin Deng, Ping Luo, Chen Change Loy, and Xiaoou Tang. Pedestrian attribute recognition at far distance. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 789–792. ACM, 2014. 2

[8] Yang Fu, Yunchao Wei, Yuqian Zhou, Honghui Shi, Gao Huang, Xinchao Wang, Zhiqiang Yao, and Thomas Huang. Horizontal pyramid matching for person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8295–8302, 2019. 2

[9] Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and Junjie Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. In *IEEE ICCV*, 2019. 2

[10] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 3

[11] Kai Han, Jianyuan Guo, Chao Zhang, and Mingjian Zhu. Attribute-aware attention model for fine-grained representation learning. In *2018 ACM Multimedia Conference on Multimedia Conference*, pages 2040–2048. ACM, 2018. 2

[12] Keke He, Yanwei Fu, Wuhao Zhang, Chengjie Wang, Yu-Gang Jiang, Feiyue Huang, and Xiangyang Xue. Harnessing synthesized abstraction images to improve facial attribute recognition. In *IJCAI*, pages 733–740, 2018. 2, 5, 8

[13] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Learning deep representation for imbalanced classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5375–5384, 2016. 3

[14] Chen Huang, Yining Li, Change Loy Chen, and Xiaoou Tang. Deep imbalanced learning for face recognition and attribute prediction. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 8

[15] Emad Sami Jaha and Mark S Nixon. Soft biometrics for subject identification using clothing attributes. In *IEEE International Joint Conference on Biometrics*, pages 1–6. IEEE, 2014. 1

[16] Mahdi M Kalayeh, Boqing Gong, and Mubarak Shah. Improving facial attribute prediction using semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6942–6950, 2017. 8

[17] Furqan M Khan and Francois Bremond. Unsupervised data association for metric learning in the context of multi-shot person re-identification. In *2016 13th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 256–262. IEEE, 2016. 3

[18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[19] Brian Kulis et al. Metric learning: A survey. *Foundations and Trends® in Machine Learning*, 5(4):287–364, 2013. 3

[20] Neeraj Kumar, Peter Belhumeur, and Shree Nayar. Facetracer: A search engine for large collections of images with faces. In *European conference on computer vision*, pages 340–353. Springer, 2008. 8

[21] Neeraj Kumar, Alexander C Berg, Peter N Belhumeur, and Shree K Nayar. Attribute and simile classifiers for face verification. In *2009 IEEE 12th International Conference on Computer Vision*, pages 365–372. IEEE, 2009. 1

[22] Ryan Layne, Timothy M Hospedales, Shaogang Gong, and Q Mary. Person re-identification by attributes. In *Bmvc*, volume 2, page 8, 2012. 1

[23] Dangwei Li, Xiaotang Chen, and Kaiqi Huang. Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. In *ACPR*, pages 111–115, 2015. 2, 5, 6

[24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 2

[25] Yutian Lin, Liang Zheng, Zhedong Zheng, Yu Wu, Zhilan Hu, Chenggang Yan, and Yi Yang. Improving person re-identification by attribute and identity learning. *Pattern Recognition*, 2019. 2, 3, 5, 6

[26] Chunxiao Liu, Shaogang Gong, Chen Change Loy, and Xinggang Lin. Person re-identification: What features are important? In *European Conference on Computer Vision*, pages 391–401. Springer, 2012. 1

[27] Hao Liu, Jingjing Wu, Jianguo Jiang, Meibin Qi, and Ren Bo. Sequence-based person attribute recognition with joint ctc-attention model. *arXiv preprint arXiv:1811.08115*, 2018. 3, 6

[28] Pengze Liu, Xihui Liu, Junjie Yan, and Jing Shao. Localization guided learning for pedestrian attribute recognition. *arXiv preprint arXiv:1808.09102*, 2018. 1, 2

[29] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 212–220, 2017. 3

[30] Xiao Liu, Jiang Wang, Shilei Wen, Errui Ding, and Yuanqing Lin. Localizing by describing: Attribute-guided attention localization for fine-grained recognition. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017. 2

[31] Xihui Liu, Haiyu Zhao, Maoqing Tian, Lu Sheng, Jing Shao, Shuai Yi, Junjie Yan, and Xiaogang Wang. Hydraplus-net: Attentive deep features for pedestrian analysis. In *Proceedings of the IEEE international conference on computer vision*, pages 350–359, 2017. 2, 6

[32] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 2, 5, 8

[33] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 5

[34] Hieu V Nguyen and Li Bai. Cosine similarity metric learning for face verification. In *Asian conference on computer vision*, pages 709–720. Springer, 2010. 1

[35] Seyoung Park, Bruce Xiaohan Nie, and Song-Chun Zhu. Attribute and-or grammar for joint parsing of human pose, parts and attributes. *IEEE transactions on pattern analysis and machine intelligence*, 40(7):1555–1569, 2017. 1

[36] Peixi Peng, Yonghong Tian, Tao Xiang, Yaowei Wang, and Tiejun Huang. Joint learning of semantic and latent attributes. In *European Conference on Computer Vision*, pages 336–353. Springer, 2016. 1

[37] Haotong Qin, Ruihao Gong, Xianglong Liu, Mingzhu Shen, Ziran Wei, Fengwei Yu, and Jingkuan Song. Ir-net: Forward and backward information retention for highly accurate binary neural networks. In *IEEE CVPR*, June 2020. 2

[38] Ethan M Rudd, Manuel Günther, and Terrance E Boult. Moon: A mixed objective optimization network for the recognition of facial attributes. In *European Conference on Computer Vision*, pages 19–35. Springer, 2016. 2, 8

[39] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Deep imbalanced attribute classification using visual attention aggregation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 680–697, 2018. 2, 6

[40] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 3, 4

[41] Patrick Sudowe, Hannah Spitzer, and Bastian Leibe. Person attribute recognition with a jointly-trained holistic cnn model. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 87–95, 2015. 2

[42] Chenxin Sun, Na Jiang, Lei Zhang, Yuehua Wang, Wei Wu, and Zhong Zhou. Unified framework for joint attribute classification and person re-identification. In *International Conference on Artificial Neural Networks*, pages 637–647. Springer, 2018. 3, 6

[43] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Advances in neural information processing systems*, pages 1988–1996, 2014. 5

[44] Yi Sun, Ding Liang, Xiaogang Wang, and Xiaoou Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015. 3

[45] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Hybrid deep learning for face verification. In *Proceedings of the IEEE international conference on computer vision*, pages 1489–1496, 2013. 1

[46] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708, 2014. 1

[47] Chufeng Tang, Lu Sheng, Zhaoxiang Zhang, and Xiaolin Hu. Improving pedestrian attribute recognition with weakly-supervised multi-scale attribute-specific localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4997–5006, 2019. 2

[48] Guangcong Wang, Jianhuang Lai, Peigen Huang, and Xiaohua Xie. Spatial-temporal person re-identification. pages 8933–8940, 2019. 2

[49] Jing Wang, Yu Cheng, and Rogerio Schmidt Feris. Walk and learn: Facial attribute representation learning from egocentric video and contextual data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2295–2304, 2016. 8

[50] Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1393, 2014. 3

[51] Yiru Wang, Weihao Gan, Wei Wu, and Junjie Yan. Dynamic curriculum learning for imbalanced data classification. *Proceedings of the IEEE international conference on computer vision*, 2019. 5

[52] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016. 3

[53] Yudong Wu, Yichao Wu, Ruihao Gong, Yuanhao Lv, Ken Chen, Ding Liang, Xiaolin Hu, Xianglong Liu, and Junjie Yan. Rotation consistent margin loss for efficient low-bit face recognition. In *IEEE CVPR*, June 2020. 3

[54] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Deep metric learning for person re-identification. In *2014 22nd International Conference on Pattern Recognition*, pages 34–39. IEEE, 2014. 3

[55] Xin Zhao, Liufang Sang, Guiguang Ding, Yuchen Guo, and Xiaoming Jin. Grouping attribute recognition for pedestrian with joint recurrent learning. In *IJCAI*, pages 3177–3183, 2018. 2

[56] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*, pages 1116–1124, 2015. 5

[57] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3754–3762, 2017. 5