

Homocentric Hypersphere Feature Embedding for Person Re-identification

Wangmeng Xiang, Jianqiang Huang, Xianbiao Qi, Xiansheng Hua, *Fellow, IEEE* and Lei Zhang, *Fellow, IEEE*

Abstract—Person re-identification (Person ReID) is a challenging task due to the large variations in camera viewpoint, lighting, resolution, and human pose. Recently, with the advancement of deep learning technologies, the performance of Person ReID has been improved swiftly. Feature extraction and feature matching are two crucial components in the training and deployment stages of Person ReID. However, many existing Person ReID methods have measure inconsistency between the training stage and the deployment stage, and they couple magnitude and orientation information of feature vectors in feature representation. Meanwhile, traditional triplet loss methods focus on samples within a mini-batch and lack knowledge of global feature distribution. To address these issues, we propose a novel homocentric hypersphere embedding scheme to decouple magnitude and orientation information for both feature and weight vectors, and reformulate classification loss and triplet loss to their angular versions and combine them into an angular discriminative loss. We evaluate our proposed method extensively on the widely used Person ReID benchmarks, including Market1501, CUHK03 and DukeMTMC-ReID. Our method demonstrates leading performance on all datasets.

Index Terms—person re-identification, deep learning, feature learning, metric learning

I. INTRODUCTION

PERSON re-identification (Person ReID) is an important computer vision task, which aims to identify a person from a set of gallery images captured under different cameras, or different timestamps under a single camera [1]. In recent years, Person ReID has drawn a lot of attentions from both academia and industry, due to its huge potential applications, such as suspect searching and multi-camera person tracking in a large-scale surveillance system. However, the task of Person ReID is extremely challenging due to the large variations in camera viewpoint, lighting, resolution, and human pose.

The key issue of Person ReID problem is how to match the query image/video with gallery images/videos. Generally speaking, both query and gallery images/videos are represented by feature vectors, which are extracted by either feature learning methods or hand crafted heuristic algorithms. This similarity of query feature vector and gallery feature vectors is then calculated by certain distance measures and the returned matching list is determined by the feature distances. Before the bloom of Convolutional Neural Networks

(CNN), various heuristic representations have been used for person re-identification, such as the local maximal occurrence representation (LOMO) [2], hierarchical Gaussian descriptor (GOG) [3]. These representations are designed to handle light variance, pose/view changes, and so on. Other works focused on similarity/metric learning techniques, which learn robust metrics under various conditions. However, as these methods using handcraft features and metrics, they are outperformed by CNN based methods.

With the recent development of Convolutional Neural Networks (CNN) [4], the performance of Person ReID has been increased dramatically. Fast evolvement in CNN architectures, such as AlexNet [5], VGGNet [6], GoogleNet [7] and ResNet [8], speeds up the development of Person ReID algorithms. Meanwhile, the increasing scale of Person ReID datasets [9], [10], [11], [12] also facilitates the study of Person ReID. Existing CNN based methods can be roughly grouped into three categories: 1) transferring and improving powerful CNN architectures to Person ReID [13], [14], [15], [16], [17], [18], where off-the-shelf feature extractors are used as parts of the network architecture; 2) designing more effective metrics [19], [20], [2], [21], [22], [23]; 3) combining priori into network architecture for fine-grained feature learning [15], [24], [25], [26], [27].

CNN with a triplet loss or classification loss is a popular framework for Person ReID. Many state-of-the-art methods [28], [29], [20], [30], [31] employ these loss functions or their variants during the training stage. For instance, Xiao *et al.* [31] trained a deep CNN from scratch using datasets from multiple sources. They employed softmax loss and domain guided dropout for training and achieved competitive performance. A multi-channel part-based CNN model under the triplet framework was proposed in [29]. This model shows superior performances on several popular benchmarks. Hermans *et al.* [28] conducted extensive experiments to validate the effectiveness of triplet loss. In [30], a two branch classification loss was employed to learn discriminative local and global features.

Despite having achieved great successes, traditional triplet loss and classification loss have a few problems. Triplet loss focuses on optimizing local distance metrics between positive pairs and negative pairs, but lacks knowledge of global feature distribution. Classification loss overlooks the intra-class variance and only maximizes inter-class variance. Meanwhile, most current classification loss based works on Person ReID lose the measure consistency between the training stage and the deployment stage. In the training stage, these methods calculate the inner product between the features x and the

Wangmeng Xiang, Xianbiao Qi and Lei Zhang are with the Department of Computing, The Hong Kong Polytechnic University, Kowloon, HongKong, China e-mail: (cswxiang@comp.polyu.edu.hk; csxqi@comp.polyu.edu.hk; cslzhang@comp.polyu.edu.hk).

Jianqiang Huang and Xiansheng Hua are with Alibaba DAMO Academy, Hangzhou, China email: (jianqiang.hjq@alibaba-inc.com; xiansheng.hxs@alibaba-inc.com).

weight vector \mathbf{w} . In the deployment stage, the features are first normalized. Given two L_2 normalized features \mathbf{x}_a and \mathbf{x}_b , their similarity is generally calculated by cosine similarity $s(\mathbf{x}_a, \mathbf{x}_b) = \mathbf{x}_a^T \mathbf{x}_b = \cos \theta_{ab}$, which is determined by the angle θ_{ab} between features \mathbf{x}_a and \mathbf{x}_b . The inconsistency mixes the orientation information and magnitude information of features, which generates a gap between the training stage and the deployment stage.

On the other hand, to achieve good retrieval performance using distance measures such as L_2 distance, the entries in the feature vector should better be independent. However, the weight vectors of fully-connected layers are usually correlated after the training. This leads to correlated entries in feature vectors, and causes inferior performance of learned feature representation. Several works have explored how to reduce the feature correlation in deep neural network training. DeCov [32] encourages diverse representations by minimizing the cross-covariance of hidden activations. Some works reduce correlation in feature entries by imposing orthogonal constraints on weight vectors. PCANet [33] which is proposed for image classification is featured by cascaded principal component analysis (PCA) filters. It uses unsupervised methods to learn orthogonal filters from raw images. In [34], Xie *et al.* applied a regularizer which utilizes orthonormality among different filter banks. SVDNet [14] achieves orthogonality by taking a restraint and relaxation iteration (RRI) training scheme, which iteratively integrates the orthogonality constraint in CNN training.

In this paper, we propose a homocentric hypersphere embedding learning approach for Person ReID. In our framework, the weight vectors and the features are normalized to two homocentric hyperspheres with the same coordinate origin. This decouples the magnitude and orientation of feature vectors and ensures the training and testing measure consistency. Based on the homocentric hyperspheres, we jointly consider the triplet loss and softmax classification loss from the perspective of angle discrimination. In triplet loss, the optimization of distances between features is reformulated to the optimization of angular distances between features. In classification loss, the posterior probability distribution will depend solely on the angle between weight vectors and features. To ensure feature orthogonality, we add a regularization term in the loss function without relying on complicated iteration methods. The angular versions of triplet loss and classification loss work well under the unified settings and they complement each other. Meanwhile, explicitly formulating the angle between the feature vectors in training avoids the measure inconsistency between training and evaluation stages, and improves the generalization power of our approach.

We evaluate our approach on three widely used Person ReID benchmarks, including Market1501 [9], CUHK03 [11], DukeMTMC-ReID [12]. Our method demonstrates leading results on all the evaluation datasets. Specifically, our approach achieves 78.56% mAP and 91.28% Rank-1 accuracy on the Market1501 dataset.

II. RELATED WORKS

A. CNN based Person Re-Identification

CNN based Person ReID methods consists of two key components: *feature learning* and *metric learning*.

The network architecture is crucial for feature learning in Person ReID. Using pre-trained networks has been proved to be an effective strategy in many applications especially when the data scale is not big enough to train a deep network from scratch. Some state-of-the-art CNNs, such as AlexNet, GoogLeNet, VGGNet and ResNet, have been employed as the feature extraction module in Person ReID [13], [14], [15], [18]. For example, Chen *et al.* [13] used AlexNet as feature extractor in their deep quadruplet network. Sun *et al.* [14] used CaffeNet and ResNet as backbone networks. GoogLeNet was employed as the base network in [15], [18]. Besides base networks, some works [25], [35], [36], [24] developed customized modules to capture human body prior. Specifically, Zhao *et al.* [25] considered human body structure information in ReID pipeline. Features of different body regions are extracted separately and merged by a tree-structured fusion network. Chen *et al.* [35] developed a multi-scale network architecture with a saliency-based feature fusion. Zhou *et al.* [36] built a part-based CNN to extract discriminative and stable features for body appearance.

In person Re-ID, various loss functions have been used for learning deep embedding representations, including verification loss [19], [37], [11], [38], contrastive loss [17], triplet loss [28], [29], [20] and quadruplet loss [13]. Yi *et al.* [19] adopted a siamese network and softmax loss to determine whether two input images belong to the same person or not. Shi *et al.* [20] trained their network using triplet loss with hard positive pairs mining. Chen *et al.* [13] designed a quadruplet loss with a margin based online hard negative mining. All these loss functions attempt to learn pairwise/ternary/quaternary distance relations. Beyond direct metric learning, some works [31], [30], [39] address the ReID problem from the perspective of classification. Such kind of works learn an embedding metric in an indirect way by constructing clear classification boundaries. Typically, classification loss contains a softmax layer with embedding features as input. Classification loss provides global classification boundaries and metric learning constructs pairwise/ternary/quaternary distance relations.

The Person ReID problem shares similarity to the face verification and recognition problems. Recently, several works [40], [41] in face recognition area consider hypersphere normalization constraints. In [40], SphereFace enforces the weights of the last classification layer lie on a unit hypersphere, which can be considered as a specific Weight Normalization [42] with the length of weight vector as 1. Liu *et al.* [41] normalized face features and optimized the distance between features and feature cluster centroids within a mini-batch. Our proposed method differs from these works by considering both feature-class center distance and feature-feature distance at the embedding learning stage. In addition, we directly learn the class center vectors, without using the mini-batch average as an approximation. Under homocentric

hypersphere settings, we build a unified angular understanding for the triplet loss and the softmax classification loss. As far as we know, our work is the first to introduce normalization on both class center vectors and features and apply angular discriminative loss for Person ReID tasks.

B. Orthogonal constraints on weight matrix

The highly correlated weight vectors of fully connected layers in CNN often results in correlations among entries of feature vectors, which would reduce retrieval performance. In order to fix this problem, several works have been proposed to introduce orthogonality in neural networks [32], [33], [34], [14]. Among them, SVDNet [14] uses a SVD layer as embedding layer to project high dimension features to a lower dimension manifold. The network is trained in a restraint and relaxation iteration (RRI) scheme, which iteratively integrates the orthogonality constraint in CNN training. During the restraint stage, the weight matrix is replaced by an orthogonal matrix, which is the product of the left unitary matrix and the singular value matrix, and the rest of the network is updated with this matrix fixed. At the relaxation stage, the restraint is removed and the whole network is tuned end-to-end. This procedure is repeated several times until the convergence criteria are met. The SVD layer can be learnt by back propagation. However, the training procedures are quite complicated and the iterative training could take a long time. Our method differs from previous methods by directly adding an orthogonal regularization term to the loss function, which is structurally simple, effective, and efficient in training.

III. HOMOCENTRIC HYPERSPHERE EMBEDDING

A. Triplet Loss and Softmax Loss

We briefly discuss triplet loss and softmax loss to draw forth the proposed angular version of these loss functions under the homocentric hypersphere assumption.

Triplet Loss. Triplet loss was firstly proposed in FaceNet [43] to improve face recognition and verification. It is originally derived from the Large Margin Nearest Neighbor (LMNN) method [44]. The objective of triplet loss is to learn an embedding function $f : \mathbf{I} \mapsto \mathbf{x}$ so that the embedded features of images from the same class are closer than the embedded features of images from different classes in the embedding space. A sample triplet $(\mathbf{I}_a, \mathbf{I}_p, \mathbf{I}_n)$ consists of an anchor sample image \mathbf{I}_a , a positive sample \mathbf{I}_p and a negative sample \mathbf{I}_n . $(\mathbf{x}_a, \mathbf{x}_p, \mathbf{x}_n)$ are the corresponding embedded features of $(\mathbf{I}_a, \mathbf{I}_p, \mathbf{I}_n)$. The general formulation for the triplet loss can be represented as follows:

$$\mathcal{L}_t = \sum_{i \in B} \left[\underbrace{\|\mathbf{x}_a^i - \mathbf{x}_p^i\|_2^2}_{\text{pull}} - \underbrace{\|\mathbf{x}_a^i - \mathbf{x}_n^i\|_2^2}_{\text{push}} + m \right]_+, \quad (1)$$

where $[\sigma]_+$ denotes $\max(\sigma, 0)$, m is a preset margin, and B is the number of triplets.

Eq. 1 contains a pull term to pull samples from the same class together, and a push term to push samples from different

classes away. The training process of the neural network is to optimize the embedding function f (adjust network parameters) to ensure that, given the feature \mathbf{x}_a^i of an anchor sample \mathbf{I}_a in the i -th triplet, the distance between \mathbf{x}_p^i and \mathbf{x}_a^i is smaller than the distance between \mathbf{x}_n^i and \mathbf{x}_a^i by a margin m .

One key issue in training Deep CNN with the triplet loss is hard triplet mining. The learning process may be dominated by simple triplets, and thus fail to learn an effective embedding. In this work, we employ an online hard triplet mining strategy in which we only consider the hardest triplets within a mini-batch. During the training phase, for each mini-batch, we randomly select P identities, and for each identity, we randomly select N samples. Therefore, each mini-batch consists of $P \times N$ samples. For each sample, we construct its hard triplet by choosing its farthest positive sample feature as \mathbf{x}_p and its closest negative sample feature as \mathbf{x}_n . In this way, we modify the triplet loss with the online hard triplet mining as

$$\mathcal{L}_{ht} = \sum_{i \in P \times N} \left[\|\mathbf{x}_a^i - \mathbf{x}_{fp}^i\|_2^2 - \|\mathbf{x}_a^i - \mathbf{x}_{cn}^i\|_2^2 + m \right]_+, \quad (2)$$

where \mathbf{x}_{fp}^i denotes the farthest one among the $N - 1$ positive sample features, and \mathbf{x}_{cn}^i stands for the closest one among the $(P - 1) \times N$ negative sample features for the anchor \mathbf{x}_a^i .

From Eq. 2, we can see that the effect of triplet loss with online hard example mining is largely affected by the mini-batch size. With a larger batch size and more samples for each identity, it is more likely to find a farther positive sample and closer negative sample, resulting in a harder triplet. It is also worth noticing that triplet loss is effective in learning pairwise/ternary distance relations; however, it only considers samples within a mini-batch without considering global feature distribution. The difficulty of training deep CNN with triplet loss will increase with the growth of dataset scale, because the number of triplets will increase exponentially with the increase of the number of training samples.

Softmax Loss. Softmax loss [45] converts an input feature into a posterior probability distribution. In softmax loss, the predicted posterior probability for the ℓ -th class is calculated as follows:

$$p_\ell = \frac{\exp(z_\ell)}{\sum_{k=1}^K \exp(z_k)}, \quad (3)$$

with

$$z_\ell = \mathbf{w}_\ell^T \mathbf{x} + b_\ell, \quad (4)$$

where \mathbf{w}_ℓ and b_ℓ are weight vector and bias of the last fully-connected layer for the ℓ -th class, K is the total number of classes.

With the posterior probability distribution of each input feature, a cross-entropy loss can be calculated as follows:

$$\mathcal{L}_c = \sum_i -\log(p_{y_i}), \quad (5)$$

where y_i is the label of the i -th input sample.

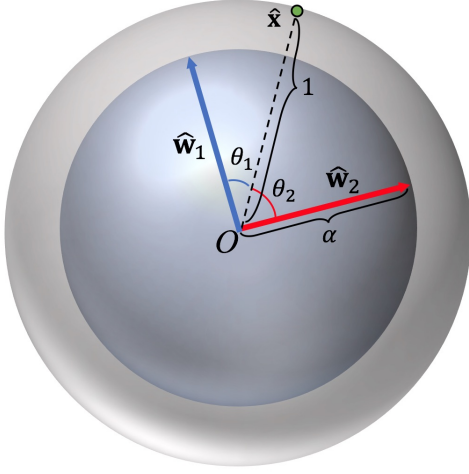


Fig. 1. Illustration of two homocentric hyperspheres. The small blue hypersphere is the weight hypersphere and the gray large hypersphere is the feature hypersphere. Weight vectors $\hat{\mathbf{w}}_\ell$ are represented by arrow lines with magnitude α and different colors denote for different classes. Feature vectors $\hat{\mathbf{x}}$ lie on the outer hypersphere with unit radius. θ_1 and θ_2 are angle distances between the feature $\hat{\mathbf{x}}$ and weight vectors $\hat{\mathbf{w}}_1$ and $\hat{\mathbf{w}}_2$, respectively.

As indicated in the introduction section, in the deployment stage, the angle distances between the features themselves, and between the features and the weight vectors are crucial. However, the original triplet loss and softmax loss do not explicitly take the angle distance into account. There is no previous work targeting to address this problem in Person ReID. To overcome this issue, we propose a novel homocentric hypersphere feature embedding learning method, as described in the next subsection.

B. Homocentric Hypersphere Constrained Embedding Learning

Homocentric Hypersphere. The angle distances between features, and between features \mathbf{x} and weight vectors \mathbf{w}_ℓ are discriminative. To ensure the features \mathbf{x} and the weight vectors \mathbf{w}_ℓ stay in a unified coordinate space, we propose a homocentric hypersphere feature embedding scheme. In our homocentric hypersphere, the weight vectors \mathbf{w}_ℓ and the features \mathbf{x} lie on two individual hyperspheres but with the same origin. The proposed homocentric hypersphere is defined as

$$\hat{\mathbf{w}}_\ell = \alpha \frac{\mathbf{w}_\ell}{\|\mathbf{w}_\ell\|}, \hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|}, \quad (6)$$

where $\hat{\mathbf{w}}_\ell$ is an L_2 normalized weight vector scaled by a factor α , and $\hat{\mathbf{x}}$ is an L_2 normalized unit feature vector.

As shown in Fig. 1, the weight vectors $\hat{\mathbf{w}}_\ell$ lie on a hypersphere with radius α , and the feature $\hat{\mathbf{x}}$ lies on a hypersphere with unit radius. These two hyperspheres have the same origin \mathbf{o} . In the hypersphere space, closer features have a smaller angle distance. A feature will be more likely assigned to the category whose weight vector has a smaller angle with the feature. As shown in Fig. 1, the feature $\hat{\mathbf{x}}$ will be more likely assigned to class 1 since it has a smaller angle distance θ_1 to $\hat{\mathbf{w}}_1$.

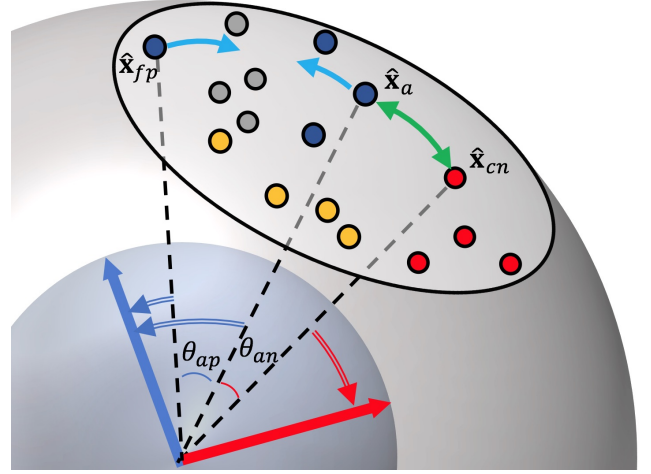


Fig. 2. The illustration of joint angular loss. The small circles within gray ellipse are features extracted from the same mini-batch. Different colors stand for different labels and all the features lie on the same feature hypersphere. Classification centers are represented by blue and red arrow lines. The green arrows show the pushing term of the triplet, while light blue arrows show the pulling term. The blue and red arrows with double lines demonstrate the guidance effect of class center vectors.

Based on the proposed homocentric hypersphere, we reformulate the traditional triplet loss and classification loss into a new angular triplet loss and a new angular classification loss, respectively. Finally, we consider these two new losses in a unified angular framework.

Angular Triplet Loss. Distance measure on triplet loss is essential for feature learning. In the homocentric hypersphere space, the features are normalized unit vectors. Given features $\hat{\mathbf{x}}_1$ and $\hat{\mathbf{x}}_2$, $\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2$ is a chord on feature hypersphere. As three edges of a triangle $(\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2)$ are known, the angle between features $(\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2)$ can be represented as

$$\theta = 2 \arcsin\left(\frac{\|\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2\|_2^2}{2}\right), \quad (7)$$

where $\theta \in [0, \pi]$. As shown in Fig. 2, given the triangle $(\hat{\mathbf{x}}_a, \hat{\mathbf{x}}_{cn}, \hat{\mathbf{x}}_a - \hat{\mathbf{x}}_{cn})$, the angle θ_{an} between $\hat{\mathbf{x}}_a$ and $\hat{\mathbf{x}}_{cn}$ can be calculated using Eq. 7.

Under the definition of homocentric hypersphere, it is more natural to measure the angle distance between features rather than Euclidean distance. Therefore, we reformulate Eq. 2 into a new angular triplet loss as

$$\mathcal{L}_{at} = \sum_{i \in P \times N} [\theta_{ap}^i - \theta_{an}^i + \theta_m]_+, \quad (8)$$

where θ_{ap}^i stands for the angle between the anchor feature and the hardest (farthest) positive feature. θ_{an}^i represents the angle between the anchor feature and the hardest (closest) negative feature. θ_m is an angular margin.

As illustrated in Fig. 2, we can see that the role of push term in Eq. 1 is now to enlarge θ_{an} and the role of pull term in Eq. 1 is to minimize θ_{ap} w.r.t. angular margin θ_m .

Angular Classification Loss. In the homocentric hypersphere formulation, the features and the class center vectors are enforced to have the same origin \mathbf{o} . To ensure the requirement,

we set the bias term b in the original softmax loss function to 0. In this way, we can reformulate Eq. 4 as

$$z_\ell = \alpha \frac{\mathbf{w}_\ell^T \mathbf{x}}{\|\mathbf{w}_\ell\| \|\mathbf{x}\|} = \alpha \cos \theta_\ell, \quad (9)$$

where θ_ℓ is the angle between \mathbf{w}_ℓ and \mathbf{x} . In this case, the original softmax loss function in Eq. 5 can be reformulated as

$$\mathcal{L}_{ac} = \sum_i -\log\left(\frac{\exp(\alpha \cos \theta_{y_i})}{\sum_{k=1}^K \exp(\alpha \cos \theta_k)}\right). \quad (10)$$

According to Eq. 9 and Eq. 10, in the training stage, our target is to maximize the cosine similarity between the features and the weight vectors of their corresponding class. This is actually minimizing the angle distance between them. The weight vector \mathbf{w}_ℓ will accumulate information from all the samples for training, and can be regarded as the class center vector of class ℓ . The posterior probability p_ℓ for class ℓ now depends solely on the angle θ_{y_i} between feature vector \mathbf{x} and class center vector \mathbf{w}_ℓ . In the testing stage, the input embedding feature will rank its neighbors according to the angle distances. Therefore, the similarity measures in the training stage and the testing stage are consistent in nature.

Joint Angular Loss. In the previous development, we have reformulated the triplet loss and classification loss into their angular versions. Here, we combine these two loss functions as

$$\mathcal{L}_a = \mathcal{L}_{at} + \lambda \mathcal{L}_{ac}, \quad (11)$$

where \mathcal{L}_{at} and \mathcal{L}_{ac} stand for angular triplet loss and angular classification loss, respectively. λ is a trade-off parameter between these two loss functions.

The proposed joint angular loss in Eq. 11 is a natural way to combine triplet loss and softmax loss under a unified angular framework. \mathcal{L}_{ac} draws clear classification boundaries between categories, and \mathcal{L}_{at} gives a specific pairwise/ternary angle relation. The angular margin θ_m in \mathcal{L}_{at} has a direct connection to feature discrimination on the person feature manifold, while the identity classification center distribution in \mathcal{L}_{ac} provides extra information to guide feature embedding learning.

In Fig. 2, we illustrate how joint angular loss works. We view the optimization of the joint angular loss from an angular discrimination perspective. For each mini-batch, we construct hard triplet by selecting the hardest positive feature $\hat{\mathbf{x}}_{fp}$ with the longest angle distance to anchor feature $\hat{\mathbf{x}}_a$, and the hardest negative feature $\hat{\mathbf{x}}_{cn}$ with the smallest angle distance to $\hat{\mathbf{x}}_a$. The angular triplet loss maximizes θ_{an} and minimizes θ_{ap} w.r.t. to angular margin θ_m . The angular classification loss minimizes the angles θ_{y_i} between features and their corresponding class center vectors. The class center vectors from angular classification loss guide the feature embedding learning. This accelerates model training by preventing optimization from getting stuck in local metric learning.

C. Orthogonal Constraints on Embedding Layer

The backbone CNN is often followed by a linear embedding layer to project high dimensional features into a low dimensional feature space. The correlation among weight vectors

in the embedding layer can compromise the performance of descriptors. It is important to reduce the redundancy of weight vectors in the embedding layer. Many works [32], [33], [34], [14] show that an orthogonal constraints on the weight matrix is effective to achieve this goal. This motivates us to place orthogonal constraints on the embedding layer by forcing $\mathbf{W}_e^T \mathbf{W}_e$ to be a diagonal matrix, where \mathbf{W}_e represents the weight of embedding layer. Specifically, we define

$$R_e = \sum \|\mathbf{W}_e^T \mathbf{W}_e - \mathbf{I}\| \quad (12)$$

and add this regularization term into the loss function. The whole loss function now becomes

$$\mathcal{L} = \mathcal{L}_a + \gamma R_e, \quad (13)$$

where γ is a trade-off parameter.

We adopt the metric used in [14] to measure the orthogonality of embedding weights:

$$S(\mathbf{W}_e) = \frac{\sum_{i=1}^k g_{ii}}{\sum_{i=1}^k \sum_{j=1}^k |g_{ij}|}, \quad (14)$$

where g_{ij} denotes the entries in $\mathbf{W}_e^T \mathbf{W}_e$ and k represents the number of weight vectors in \mathbf{W}_e . The values of $S(\mathbf{W}_e)$ range within $[\frac{1}{k}, 1]$. A higher value of $S(\mathbf{W}_e)$ indicates better orthogonality.

We conduct cross validation to find how the embedding layer's orthogonality affect the performance of trained models, and determine the parameter γ . More details can be found in section IV.D.

D. The Flowchart of Algorithm

The flowchart of our algorithm is shown in Fig. 3. We divide our model into two parts: feature extraction part and metric learning part. For feature extraction, different CNN architectures can be employed, such as SqueezeNet, VGG and ResNet. We use ResNet50 as the default feature extraction network. All the CNNs are pre-trained on the ImageNet [46] dataset. We apply average pooling on the final convolutional feature maps to obtain feature vectors. For the metric learning strategy, we employ an embedding layer after average pooling to generate low-dimensional features. The homocentric hypersphere constraints are then applied on features and class center vectors. Joint angular loss is employed for metric learning with additional orthogonal regularization.

IV. EXPERIMENTS

A. Datasets

We evaluate our approach on three widely used Person ReID benchmarks, including Market1501 [9], CUHK03 [11] and DukeMTMC-ReID [12].

Market1501 contains 32,668 pedestrian images of 1,501 identities. The images were collected by five high-resolution and one low-resolution cameras. The resolution of pedestrian images is 64×128 . Each identity in the dataset appears on at least two different views. Following the official split of training

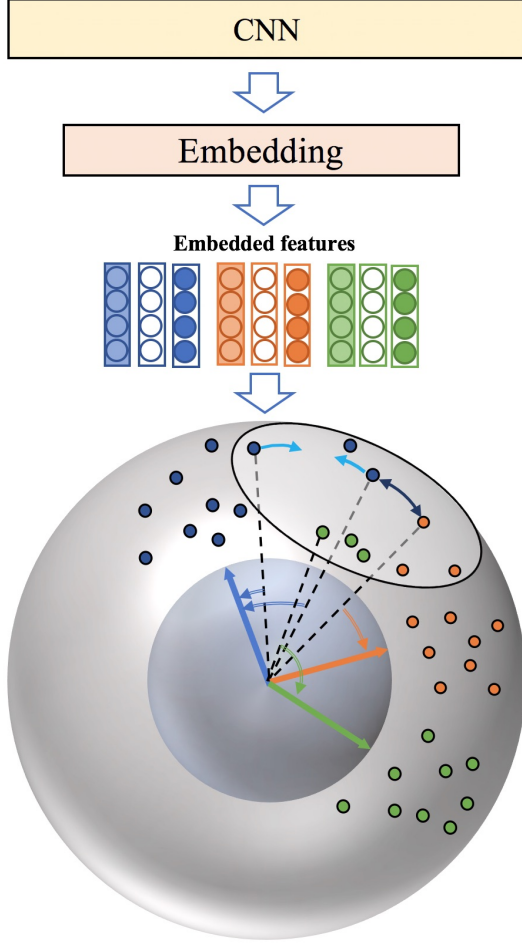


Fig. 3. The flow chart of our proposed algorithm. The top part is a CNN for feature extraction and the middle part is the embedding layer. The bottom illustrates our proposed homocentric hypersphere embedding learning part. The different colors of embedded features represent different identities. The embedded features are represented by the points on the feature hypersphere.

and testing sets, we use 751 identities for training, and the rest 750 for testing.

CUHK03 has 1,467 identities from two different views in the CUHK campus. There are 14,097 images obtained by Deformable Part based Models (DPM) [47], and 14,096 images obtained by manually labeling. 1,367 person identities are used for training, and the rest 100 identities for testing. In this paper, we follow the official evaluation protocols and report performance of our method using both DPM detected images and manually labeled images.

DukeMTMC-ReID is a subset of the multi-target, multi-camera pedestrian tracking dataset. In this work, we use a subset of [12] provided by [48], which has 16,522 images from 702 identities for training. The images were captured by eight different high-resolution cameras. The images of pedestrians are manually cropped. For testing, it has 2,228 images for querying and 17,661 gallery images from 702 identities. We follow the evaluation protocol in [48].

Table I provides a statistical summary of each dataset. It lists the number of identities (ID), bounding boxes (BBboxes),

TABLE I
THE STATISTICS OF PERSON REID DATASETS.

Datasets	ID	BBboxes	Distra	Views
Market1501	1501	32668	2793	6
CUHK03	1467	14097	0	2
DukeMTMC-ReID	702	16522	0	8

distractors (Distra), and views (Views) in each dataset. Fig. 4 shows some sample images for each dataset.

B. Evaluation Protocol

We use Single-Query (SQ) by default for all datasets. According to previous evaluation protocols used in each dataset, we adopt three evaluation protocols, including Cumulated Matching Characteristics (CMC) (top1, top5 and top10), mean Average Precision (mAP), and Rank-1 identification rate. On Market1501 and DukeMTMC-ReID, query and gallery sets may have the same camera views, but for each individual query identity, his/her gallery samples from the same cameras are excluded. On CUHK03, we follow the standard testing protocol. For each query, we randomly sample one instance for each gallery identity, and compute a CMC curve in the single-gallery-shot setting. As random selection is involved, we repeat the evaluation procedure for 10 times and report the mean results. On Market1501, we report CMC and mAP. On CUHK03, we report CMC for both detected and labeled datasets. On DukeMTMC-ReID, we report Rank-1 identification rate and mAP.

C. Implementation Details

Data preprocessing. We use the same data pre-processing methods on all datasets. In the training stage, common data augmentation methods are applied to images, including random flipping, shifting, zooming, cropping and random erasing [49]. The images are then resized to 256×128 . As we use pre-trained networks from torchvision¹, all images' pixel values are normalized to $[0, 1]$, subtracted by mean pixel values of RGB channels and then divided by standard deviation of each channel.

Optimization. For joint angular loss, our model is trained with the batch size equals 256, and for each instance we randomly select 8 samples. We apply Adam optimizer and set the original learning rate to 1×10^{-3} for the first 50 epochs and gradually decrease it to 1×10^{-4} for next 50 epochs and 10^{-5} for the last 50 epochs. For homocentric hypersphere embedding learning with orthogonal constraints, extensive experiments are conducted to illustrate the effect of hyper parameters of the model. The detailed analysis and discussion on parameter setting would be given in the next subsection. We use the same training schedule for triplet loss baseline. For classification loss baseline, we apply SGD momentum with Nesterov. We set initial learning rate to 0.1 for metric learning layers and 0.01 for pre-trained feature

¹<https://github.com/pytorch/vision>



(a) Market1501



(b) CUHK03



(c) DukeMTMC

Fig. 4. Sample images from Market1501, CUHK03 and DukeMTMC datasets.

extraction layer. We decrease the learning rate by a factor of 10 for every 40 epochs until convergence. Our implementation is built on modified Open-ReID², an open source Person ReID library.

In the following subsections, we abbreviate our method with the proposed Joint Angular Loss as JAL.

D. Exploratory Experiments

In this section, we conduct experiments to determine the values of hyperparameters in our model, including the trade-off parameter λ between angular triplet loss and classification loss, angular margin θ_m , scaling parameters α and orthogonal constraint parameter γ . We further show the influences of different loss functions, orthogonal constraints and test data augmentation to the ReID performance.

Choice of λ and angular margin θ_m . As we can see from Eq. 8 and Eq. 11, λ controls the trade-off between angular triplet loss and classification loss, while angular margin θ_m controls the boundary distance of identities on feature embedding manifold. We cross validate λ and θ_m on the ReID results on Market1501. The value of λ is selected from

TABLE II
INFLUENCES OF DIFFERENT λ AND θ_m ON MODEL PERFORMANCE ON VALIDATION SET OF MARKET1501.

λ	θ_m	Top1	Top5	Top10	mAP
0.1	1	89.70	96.47	97.95	75.18
0.1	3	90.11	96.73	98.01	76.10
0.1	5	89.82	96.59	97.98	76.16
0.1	8	90.05	96.29	97.45	75.96
0.1	10	89.79	96.47	97.95	76.22
0.2	1	91.11	96.52	98.00	77.88
0.2	3	90.95	96.78	97.94	78.05
0.2	5	91.01	96.51	97.92	77.77
0.2	8	90.53	96.38	97.73	77.31
0.2	10	90.65	96.14	97.63	77.13
0.5	1	90.83	96.59	97.71	77.36
0.5	3	90.97	96.50	97.92	77.76
0.5	5	91.06	97.12	98.04	77.81
0.5	8	90.56	96.44	97.83	77.78
0.5	10	90.47	96.18	97.77	76.40
1.0	1	90.38	96.23	97.77	76.48
1.0	3	90.68	96.53	97.92	77.22
1.0	5	90.62	96.35	97.71	76.08
1.0	8	89.55	95.87	97.71	75.23
1.0	10	90.08	96.02	97.51	76.82

TABLE III
INFLUENCES OF DIFFERENT α ON MODEL PERFORMANCE ON MARKET1501.

α	Top1	Top5	Top10	mAP
3	85.81	93.68	96.02	69.20
6	87.00	94.54	96.56	73.89
9	89.80	95.84	97.28	76.04
12	90.95	96.78	97.94	78.05
15	90.94	96.94	97.88	77.24
18	90.52	96.82	98.03	76.96
20	90.80	96.74	98.10	76.53
30	89.64	96.37	97.69	75.36

$\{0.1, 0.2, 0.5, 1\}$ and θ_m is selected from $\{1^\circ, 3^\circ, 5^\circ, 8^\circ, 10^\circ\}$. We experimentally found that models with reasonably higher weight of triplet loss and smaller angular margin tend to have better performance on validation set. As we can see from Tab. II, the best λ and θ_m on the Market1501 are around 0.2 and 3° , respectively. The Table also shows that our model is insensitive to these two parameters. To simplify the experiment, we set $\lambda = 0.2$ and $\theta_m = 3^\circ$ in all the later experiments.

Feature and weights scaling. According to Eq. 10, there is a hyper parameters α to adjust the scale of the inputs to the angular softmax layer after the feature vectors and weights normalization. We conducted experiments on different α , to study how the scaling parameter affect the performance. As we apply angular triplet loss, the value of α would not affect the triplet part. We select the value of α from $\{3, 6, 9, 12, 15, 18, 20, 30\}$. As it is showed in Tab. III, the

²<https://cysu.github.io/open-reid/>

TABLE IV
INFLUENCES OF DIFFERENT γ ON MODEL PERFORMANCE ON MARKET1501. SW IS THE MEASURE OF ORTHOGONALITY OF WEIGHTS.

γ	Top1	Top5	Top10	mAP	SW
10^{-1}	91.02	96.60	97.98	77.92	0.9995
10^{-2}	91.30	96.52	98.10	78.37	0.9995
10^{-3}	91.30	96.60	97.70	78.50	0.9986
10^{-4}	91.07	96.46	98.06	78.26	0.9984
0	90.95	96.78	97.94	78.05	0.0816

model performance is much influenced by the value of α . However, the trained models generally work well in a range of α . For simplicity, we keep $\alpha = 12$ in all our experiments, and our proposed method works very robustly.

Orthogonal constraints. According to Eq. 13, the hyper parameters γ controls the trade-off between the joint angular loss and regularization term. We conducted extensive experiments to show γ 's influence on embedding layer's orthogonality. The value of γ is selected from $\{0, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$. As we can see from Tab. IV, our model is insensitive to the value of γ from 10^{-4} to 10^{-1} , while the introduction of orthogonality regularization improves the model performance (e.g., comparing $\gamma = 10^{-3}$ with $\gamma = 0$).

Effect of joint angular loss. To show the effect of different components of JAL, we conduct experiments on different variants of JAL on the Market1501, CUHK03 labeled and DukeMTMC-ReID datasets. The results are reported in Tab. V, where C stands for JAL with only classification loss. T stands for JAL with only triplet loss and hard example mining. C+T stands for JAL with combined classification loss and triplet loss. JAL represents joint angular loss without orthogonal regularization, and JAL_o represents JAL with orthogonal regularization.

Several important observations could be made from Tab. V. 1) Combining classification loss and triplet loss largely improves the performance over using them individually, with around 7.7% and 3.0% increases on mAP on Market1501, 10.2% and 1.0% increases on Top1 on CUHK03, 9.6% and 1.2% increases on mAP on DukeMTMC-reID, respectively. This shows that these two loss functions are complementary in nature. 2) Using the proposed homocentric hypersphere embedding, JAL outperforms the baseline C+T by 4% on mAP on Market1501, 0.4% on Top1 on CUHK03, and 3% on mAP on DukeMTMC. This demonstrates the benefit of Joint Angular Loss, which optimizes angle distances rather than Euclidean distances. 3) The orthogonal regularization on weights provides further performance boosts from 0.5% to 2% on CMC Top1 and mAP on the three datasets.

Test data augmentation. Test data augmentation simulates different viewpoints and occlusion effect of original person image. We apply common data augmentation methods, such as random cropping and flipping. The final feature vector of a given image is produced by averaging all features generated by the augmented images and the original one. In Tab. VI, we show the number of augmented images used for producing the final feature and the ReID performances on Market1501. One

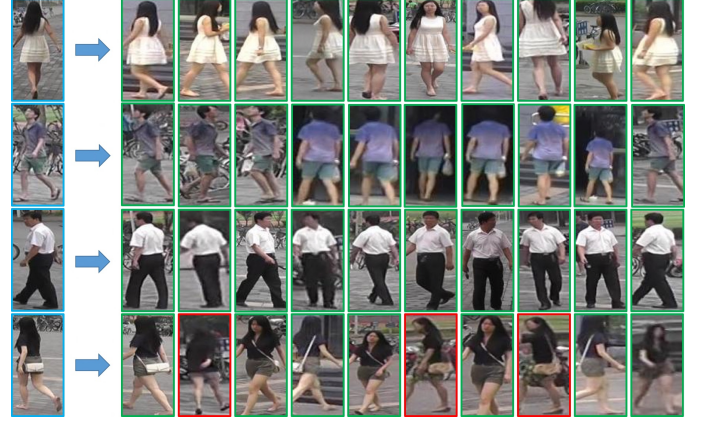


Fig. 5. Sample retrieval results on Market1501 using the proposed method. The images in the first column are the query images. The top-10 retrieved images are sorted according to the similarity scores from left to right. Gallery images captured from the same camera view as query images are already excluded from the ranking list. The correct matches are in the green rectangles, and the false matching images are in the red rectangles.

can see that data augmentation can improve the performance. In the experiments, we report the results of our method with and without data augmentation for more comprehensive comparison with other methods.

E. Comparison with state-of-the-art methods

Market1501. Market1501 is currently the largest benchmark dataset for Person ReID, and many methods have been reported on this dataset. We compare the proposed method with most of the state-of-the-arts, including Discriminative Null Space (DNS) [50], Gated Siamese Convolutional Neural Network (G-CNN) [16], Unlabeled Sample Generation GAN (GAN) [48], Deep Transfer Learning (DTL) [18], Joint Learning Multi-Loss (JLML) [30], TriNet [28], Deep Context-aware Features (DCF) [24], Spindle network (Spindle) [25], Supervised Smooth Manifold (SSM) [51], Point Set Similarity Feature (PSSF) [52], Deeply Learned Part-Aligned representation (DLPA) [15], and Pose-driven Deep Convolutional (PDC) model [53]. The experimental results are shown in Tab. VII.

With ResNet50 as the pre-trained network, the proposed JAL approach achieves 78.1% mAP and 91.0% CMC top1. It is superior to all compared methods. Specifically, JAL outperforms TriNet [28] by 6.1% on CMC top1 and 9% on mAP. In addition, JAL also outperforms the newly proposed methods, PLPA and PDC, by a large margin. With the orthogonal regularization and 16 times test augmentation, the performance of JAL is further improved.

In Fig. 5, we show the CMC top 10 retrieval results of four query images in Market1501. We can see that JAL exhibits strong robustness to scale, viewpoint and pose. The false matchings marked in red bounding box look very similar to the query image in visual appearance and human pose. These false matchings are even very challenging for human.

CUHK03. On CUHK03, we follow the standard protocol and conduct experiments using both labeled and detected datasets. We compare our model with Filter Pair Neural

TABLE V

THE EFFECT OF DIFFERENT COMPONENTS OF JAL. C STANDS FOR JAL WITH ONLY CLASSIFICATION LOSS. T STANDS FOR JAL WITH ONLY TRIPLET LOSS AND HARD EXAMPLE MINING. C+T STANDS FOR JAL WITH COMBINED CLASSIFICATION LOSS AND TRIPLET LOSS USING $\lambda = 0.2$. JAL REPRESENTS JOINT ANGULAR LOSS WITHOUT ORTHOGONAL REGULARIZATION, AND JAL_o REPRESENTS JAL WITH ORTHOGONAL REGULARIZATION.

		Market1501				CUHK03 label			DukeMTMC			
Method	Backbone	Top1	Top5	Top10	mAP	Top1	Top5	Top10	Top1	Top5	Top10	mAP
C	ResNet50	85.18	94.30	96.11	66.36	76.83	93.26	96.37	71.97	83.95	87.59	52.93
T		86.20	94.80	96.71	71.02	86.06	97.54	98.81	76.30	87.75	91.52	61.35
C+T		88.40	95.55	96.91	74.03	86.99	97.74	99.01	77.42	88.47	91.43	62.53
JAL		90.95	96.78	97.94	78.05	87.43	97.74	98.74	80.61	90.56	93.31	65.55
JAL_o		91.28	96.55	97.81	78.56	88.80	98.20	99.50	82.54	91.16	93.76	66.85

TABLE VI

TEST DATA AUGMENTATION INFLUENCES ON MARKET1501. THE MODEL IS TRAINED WITH JOINT ANGULAR LOSS AND RESNET50 AS FEATURE EXTRACTOR.

Image Num	Top1	Top5	Top10	mAP
Original	91.28	96.55	97.81	78.56
2	91.75	96.82	98.10	79.50
4	92.40	97.10	98.37	79.90
8	92.25	97.00	98.22	80.18
12	92.34	96.91	98.22	80.22
16	92.10	97.00	98.37	80.31
32	92.16	97.18	98.34	80.36

TABLE VII

SINGLE-SHOT PERFORMANCE COMPARISON OF DIFFERENT METHODS ON MARKET1501. WE HIGHLIGHT TOP-5 RESULTS ACCORDING TO MAP.

Method	Top1	Top5	Top10	mAP
DNS [50] (CVPR16)	55.4	-	-	29.9
G-CNN [16] (ECCV16)	65.9	-	-	39.6
GAN [48] (ICCV17)	78.1	-	-	56.2
DTL [18] (Arxiv16)	83.7	-	-	65.5
JLML [30] (IJCAI17)	85.1	-	-	65.5
DCF [24] (CVPR17)	80.3	-	-	57.5
Spindle [25] (CVPR17)	76.9	91.5	94.6	-
PSSF [52] (CVPR17)	70.7	90.5	-	-
DLPA [15] (ICCV17)	81.0	92.0	94.7	63.4
PDC [53] (ICCV17)	84.1	92.7	94.9	63.4
SSM [51] (CVPR17)	82.2	-	-	68.8
TriNet [28] (Arxiv17)	84.9	94.2	-	69.1
JAL	91.0	96.8	97.9	78.1
JAL_o	91.3	96.6	97.8	78.6
JAL_o +aug16	92.1	97.0	98.4	80.3

Networks (FPNN) [11], Improved Deep Learning Architecture (IDLA) [37], Local Maximal Occurrence Representation (LOMO), Sample Specific SVM (SS-SVM) [54], Discriminative Null Space (DNS) [50], Deep Context-aware Features (DCF), Quadriplet loss (Quadruplet) [13], Spindle Network (Spindle), Supervised Smooth Manifold (SSM) [51], Multi-scale Deep Architecture (MuDeep) [55], Deeply Learned Partially aligned (DLPA) [15] and Pose-driven CNN (PDC) [53]. The



Fig. 6. Sample retrieval results on CUHK03 labeled dataset using the proposed method. The images in the first column are the query images. The top-10 retrieved images are sorted according to the similarity scores from left to right. Gallery images captured from the same camera view as query images are already excluded from the ranking list. The correct matches are in the green rectangles, and the other retrieved images are in the gray rectangles.

results are shown in Tab. VIII and Tab. IX. Fig. 6 shows some sample query results on the CUHK03 labeled dataset. The proposed method retrieves all the images of the same persons from different views and the other retrieved person images are also visually similar with the query samples.

Using the same hyper parameters as used on the Market1501, our method JAL_o achieves 88.8% CMC top1 on CUHK03 labeled dataset and 86.9% CMC top1 on CUHK03 detected dataset, which are leading results on CUHK03 dataset. It is worth noticing that DLPA [15] uses both CUHK03 and Market1501 for training. Spindle [25] uses seven data sets for training, including CUHK01, CUHK03, Market1501 and etc. PDC [53] fuses two Google Inception sub-networks to consider global feature maps and part feature maps. Our approach is trained on CUHK03 with only a single ResNet50.

DukeMTMC-ReID. On this dataset, we compare our method with Scalable Person Re-identification (BoW+KISSME) [9], Local Maximal Occurrence Representation (LOMO) [2], Improved Attribute Person ReID (APR) [56], Unlabeled Sample Generation GAN (GAN), Pedestrian Alignment Network (PAN) [26], SVDNet [14] and Deep Learning Multi-Scale Representations (DPFL) [35],

TABLE VIII
SINGLE SHOT PERFORMANCE COMPARISON OF DIFFERENT METHODS ON CUHK03 LABELED DATASET. WE HIGHLIGHT THE TOP-5 PERFORMERS ACCORDING TO CMC TOP1.

Method	Top1	Top5	Top10
FPNN [11] (CVPR14)	20.7	51.5	66.5
IDLA [37] (CVPR15)	54.7	86.5	93.9
LOMO [2] (CVPR15)	52.2	82.2	92.1
SS-SVM [54] (CVPR16)	57.0	84.8	92.5
DNS [50] (CVPR16)	58.9	85.6	92.5
DCF [24] (CVPR17)	74.2	94.3	97.5
Quadruplet [13] (CVPR17)	75.5	95.2	99.2
SSM [51] (CVPR17)	76.6	94.6	98.0
MuDeep [55] (ICCV17)	76.9	96.1	98.4
DLPA [15] (ICCV17)	85.4	97.6	99.4
Spindle [25] (CVPR17)	88.5	97.8	98.6
PDC [53] (ICCV17)	88.7	98.6	99.2
JAL	87.4	97.7	98.7
JAL _o	88.8	98.2	99.5
JAL _o +aug16	89.5	97.8	99.1

TABLE IX
SINGLE SHOT PERFORMANCE COMPARISON OF DIFFERENT METHODS ON CUHK03 DETECTED DATASET. WE HIGHLIGHT THE TOP-5 PERFORMERS ACCORDING TO CMC TOP1.

Method	Top1	Top5	Top10
FPNN [11] (CVPR14)	19.9	50.0	64.0
IDLA [37] (CVPR15)	45.0	76.0	83.5
LOMO [2] (CVPR15)	46.3	78.9	88.6
SS-SVM [54] (CVPR16)	51.2	81.5	89.9
DNS [50] (CVPR16)	54.7	84.8	94.8
DCF [24] (CVPR17)	68.0	91.0	95.4
SSM [51] (CVPR17)	72.7	92.4	96.1
MuDeep [55] (ICCV17)	75.6	94.4	97.5
PDC [53] (ICCV17)	78.3	94.8	97.2
DLPA [15] (ICCV17)	81.6	97.3	98.4
JAL	84.8	97.0	97.9
JAL _o	86.9	97.4	98.5
JAL _o +aug16	88.4	97.3	98.4

which are all state-of-the-art methods we can find in literature.

As shown in Tab. X, our model achieves much better performance than other methods. The proposed JAL_o obtains 82.5% Rank-1 accuracy and 66.9% mAP, respectively. The DPFL [35] method fuses multi-scale information by combining multiple sub-networks in the training stage. In contrast, our model is trained only on DukeMTMC-ReID with a single ResNet50. The Fig. 7 shows some retrieval examples on the DukeMTMC-ReID dataset. In the second row, the person in the wrongly retrieved images has similar clothing to the person in the query image. The third row is a very challenging example, as all the persons wear similar bright green backpacks.



Fig. 7. Retrieval examples on DukeMTMC-ReID dataset using the proposed method. The images in the first column are the query images. The top-10 retrieved images are sorted according to the similarity scores from left to right. Gallery images captured from the same camera view as query images are already excluded from the ranking list. The correct matches are in the green rectangles, and the false matching images are in the red rectangles.

TABLE X
SINGLE SHOT PERFORMANCE COMPARISON OF DIFFERENT METHODS ON DUKEMTMC-REID. WE HIGHLIGHT THE TOP-5 PERFORMERS ACCORDING TO MAP.

Method	Top1	mAP
BoW+KISSME [9] (CVPR15)	25.1	12.2
LOMO [57] (CVPR15)	30.8	17.0
APR [56] (ArXiv17)	70.7	51.9
GAN [48] (ICCV17)	67.7	47.1
PAN [26] (CVPR17)	71.6	51.5
SVDNet [14] (ICCV17)	76.7	56.8
DPFL [35] (ICCV17)	79.2	60.6
JAL	80.6	65.6
JAL _o	82.5	66.9
JAL _o +aug16	83.0	67.8

V. CONCLUSION

In this paper, we proposed a homocentric hypersphere feature embedding learning approach for Person ReID task. In our homecentric hypersphere framework, the class center vectors and the features were normalized to two individual homocentric hyperspheres with the same coordinate origin. Based on the homecentric hypersphere assumption, we reformulated the classification loss and the triplet loss into their corresponding angular versions, and thus provided a natural way to jointly consider both losses to minimize the angle between intra-class features, maximize the angle between inter-class features and minimize the angle between features and their corresponding class center vectors. To reduce the redundancy in the embedding layers' weights, we placed orthogonal regularizations on embedding layer. Detailed analysis and extensive experiments were conducted on three widely used data sets to validate the effectiveness of our approach for Person ReID. The results showed that our approach achieves state-of-the-art performance on all benchmarks by using the same hyperparameters.

REFERENCES

- [1] N. Gheissari, T. B. Sebastian, and R. Hartley, "Person reidentification using spatiotemporal appearance," in *CVPR*, 2006.
- [2] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *CVPR*, 2015.
- [3] T. Matsukawa, T. Okabe, E. Suzuki, and Y. Sato, "Hierarchical gaussian descriptor for person re-identification," in *CVPR*, 2016.
- [4] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, 1995.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [6] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [7] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [9] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *CVPR*, 2015.
- [10] W. Li, R. Zhao, and X. Wang, "Human reidentification with transferred metric learning," in *ACCV*, 2012.
- [11] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *CVPR*, 2014.
- [12] E. Ristani, F. Solera, R. Zou, R. Cucchiara, and C. Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," in *ECCV workshop*, 2016.
- [13] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: a deep quadruplet network for person re-identification," in *CVPR*, 2017.
- [14] Y. Sun, L. Zheng, W. Deng, and S. Wang, "Svdnet for pedestrian retrieval," in *ICCV*, 2017.
- [15] L. Zhao, X. Li, J. Wang, and Y. Zhuang, "Deeply-learned part-aligned representations for person re-identification," in *ICCV*, 2017.
- [16] R. R. Viorio, M. Haloi, and G. Wang, "Gated siamese convolutional neural network architecture for human re-identification," in *ECCV*, 2016.
- [17] R. R. Viorio, B. Shuai, J. Lu, D. Xu, and G. Wang, "A siamese long short-term memory architecture for human re-identification," in *ECCV*, 2016.
- [18] M. Geng, Y. Wang, T. Xiang, and Y. Tian, "Deep transfer learning for person re-identification," *arXiv preprint arXiv:1611.05244*, 2016.
- [19] D. Yi, Z. Lei, and S. Z. Li, "Deep metric learning for practical person re-identification," *arXiv preprint arXiv:1407.4979*, 2014.
- [20] H. Shi, Y. Yang, X. Zhu, S. Liao, Z. Lei, W. Zheng, and S. Z. Li, "Embedding deep metric for person re-identification A study against large variations," in *ECCV*, 2016.
- [21] C. Jose and F. Fleuret, "Scalable metric learning via weighted approximate rank component analysis," in *ECCV*, 2016.
- [22] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *CVPR*, 2016.
- [23] S. Liao and S. Z. Li, "Efficient psd constrained asymmetric metric learning for person re-identification," in *ICCV*, 2015.
- [24] D. Li, X. Chen, Z. Zhang, and K. Huang, "Learning deep context-aware features over body and latent parts for person re-identification," in *CVPR*, 2017.
- [25] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, and X. Tang, "Spindle net: Person re-identification with human body region guided feature decomposition and fusion," in *CVPR*, 2017.
- [26] Z. Zheng, L. Zheng, and Y. Yang, "Pedestrian alignment network for large-scale person re-identification," in *CVPR*, 2017.
- [27] Y. Zhang, X. Li, L. Zhao, and Z. Zhang, "Semantics-aware deep correspondence structure learning for robust person re-identification," in *IJCAI*, 2016.
- [28] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [29] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng, "Person re-identification by multi-channel parts-based cnn with improved triplet loss function," in *CVPR*, 2016.
- [30] W. Li, X. Zhu, and S. Gong, "Person re-identification by deep joint learning of multi-loss classification," in *IJCAI*, 2017.
- [31] T. Xiao, H. Li, W. Ouyang, and X. Wang, "Learning deep feature representations with domain guided dropout for person re-identification," in *CVPR*, 2016.
- [32] M. Cogswell, F. Ahmed, R. B. Girshick, L. Zitnick, and D. Batra, "Reducing overfitting in deep networks by decorrelating representations," in *ICLR*, 2016.
- [33] T. H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma, "Pcanet: A simple deep learning baseline for image classification?" in *IEEE Transactions on Image Processing*, vol. 24, no. 12, Dec 2015, pp. 5017–5032.
- [34] D. Xie, J. Xiong, and S. Pu, "All you need is beyond a good init: Exploring better solution for training extremely deep convolutional neural networks with orthonormality and modulation," in *CVPR*, 2017.
- [35] Y. Chen, X. Zhu, and S. Gong, "Person re-identification by deep learning multi-scale representations," in *CVPR*, 2017.
- [36] S. Zhou, J. Wang, R. Shi, Q. Hou, Y. Gong, and N. Zheng, "Large margin learning in set to set similarity comparison for person re-identification," *IEEE Transactions on Multimedia*, 2017.
- [37] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in *CVPR*, 2015.
- [38] L. Wu, C. Shen, and A. v. d. Hengel, "Personnet: Person re-identification with deep convolutional neural networks," *arXiv preprint arXiv:1601.07255*, 2016.
- [39] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification: Past, present and future," *arXiv preprint arXiv:1610.02984*, 2016.
- [40] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *CVPR*, 2017.
- [41] Y. Liu, H. Li, and X. Wang, "Rethinking feature discrimination and polymerization for large-scale recognition," in *NIPS 2017 workshop*, 2017.
- [42] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *NIPS*, 2016.
- [43] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *CVPR*, 2015.
- [44] K. Q. Weinberger, J. Blitzer, and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," in *NIPS*, 2006.
- [45] C. M. Bishop, *Pattern recognition and machine learning*, 2006.
- [46] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009.
- [47] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *TPAMI*, 2010.
- [48] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *CVPR*, 2017.
- [49] G. K. S. L. Y. Y. Zhun Zhong, Liang Zheng, "Random erasing data augmentation," *arXiv preprint arXiv:1708.04896*, 2017.
- [50] L. Zhang, T. Xiang, and S. Gong, "Learning a discriminative null space for person re-identification," in *CVPR*, 2016.
- [51] S. Bai, X. Bai, and Q. Tian, "Scalable person re-identification on supervised smoothed manifold," 2017.
- [52] S. Zhou, J. Wang, J. Wang, Y. Gong, and N. Zheng, "Point to set similarity based deep feature learning for person re-identification," in *CVPR*, 2017.
- [53] C. Su, J. Li, S. Zhang, J. Xing, W. Gao, and Q. Tian, "Pose-driven deep convolutional model for person re-identification," in *CVPR*, 2017.
- [54] Y. Zhang, B. Li, H. Lu, A. Irie, and X. Ruan, "Sample-specific svm learning for person re-identification," in *CVPR*, 2016.
- [55] X. Qian, Y. Fu, Y.-G. Jiang, T. Xiang, and X. Xue, "Multi-scale deep learning architectures for person re-identification," in *ICCV*, 2017.
- [56] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, and Y. Yang, "Improving person re-identification by attribute and identity learning," *arXiv preprint arXiv:1703.07220*, 2017.
- [57] X. Z. Shengcai Liao, Yang Hu and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *CVPR*, 2015.