# Deep Ranking for Person Re-identification via Joint Representation Learning

Shi-Zhe Chen, Chun-Chao Guo, *Student Member, IEEE,* and Jian-Huang Lai, *Senior Member, IEEE*

*Abstract*—This paper proposes a novel approach to person re-identification, a fundamental task in distributed multi-camera surveillance systems. Although a variety of powerful algorithms have been presented in the past few years, most of them usually focus on designing hand-crafted features and learning metrics either individually or sequentially. Different from previous works, we formulate a unified deep ranking framework that jointly tackles both of these key components to maximize their strengths. We start from the principle that the correct match of the probe image should be positioned in the top rank within the whole gallery set. An effective learning-to-rank algorithm is proposed to minimize the cost corresponding to the ranking disorders of the gallery. The ranking model is solved with a deep convolutional neural network (CNN) that builds the relation between input image pairs and their similarity scores through joint representation learning directly from raw image pixels. The proposed framework allows us to get rid of feature engineering and does not rely on any assumption. An extensive comparative evaluation is given, demonstrating that our approach significantly outperforms all state-of-the-art approaches, including both traditional and CNN-based methods on the challenging VIPeR, CUHK-01 and CAVIAR4REID datasets. Additionally, our approach has better ability to generalize across datasets without fine-tuning.

*Index Terms*—Person re-identification, deep convolutional neural network, learning to rank.

## I. INTRODUCTION

**P**ERSON re-identification underpins many critical applications in long-term multi-camera tracking [1] and forensic search [2], and is increasingly receiving attention as a key component of video surveillance [3]. Given an image of a target pedestrian captured by one camera, a person re-identification system attempts to recognize the occurrence of that target from a gallery of already-labeled subjects. Since the camera views of the realistic video surveillance system are usually disjoint, the system has to re-identify pedestrians based solely on visual cues most of the time. However, the appearance of a given individual undergoes drastic changes owing to complex variations in illumination, pose, viewpoint,

S.-Z. Chen is with the School of Electronics and Information Technology, Sun Yat-sen University, Guangzhou 510006, China, and also with the Guangdong Province Key Laboratory of Information Security, China (e-mail: chenshizh@gmail.com).

C.-C. Guo and J.-H. Lai are with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou 510006, China. C.-C. Guo is also with the Guangdong Provincial Key Laboratory of Digital Signal and Image Processing Techniques, China. J.-H. Lai is also with the Guangdong Province Key Laboratory of Information Security, China (e-mail: chunchaoguo@gmail.com; stsljh@mail.sysu.edu.cn).
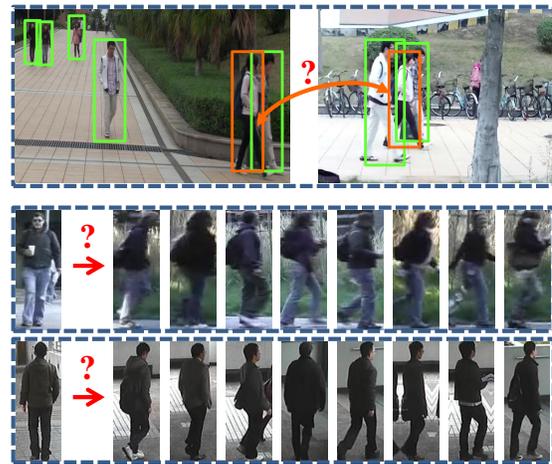


Fig. 1. Person re-identification remains challenging due to the drastic cross-view variations caused by illumination, occlusion, pose, etc. The images in the first row are taken from realistic surveillance systems, and those in the second and third rows come from the VIPeR and CUHK-01 datasets, respectively.

occlusion, image resolution and camera setting, rendering person re-identification an unsolved and challenging problem (Figure 1).

Because the main difficulty in person re-identification arises from severe changes across non-overlapping camera views, an obvious solution is to design robust and discriminative descriptors for cross-view matching. Low-level features such as color (color histograms of different color spaces [4]–[7]) and texture (LBP [8], Gabor [6], [8], [9]) are commonly used for this purpose. Some studies have sought more distinct and reliable feature representation for pedestrians, including symmetry-driven accumulation [10], horizontal stripe-based partition [6], [11], pyramid matching [12], and salience matching [13], [14]. Unfortunately, it is extremely difficult to design a feature that is distinct, reliable and invariant to severe changes and misalignment across disjoint views.

Person re-identification has also been cast as a metric learning problem, resulting in significant performance improvements [4]–[6], [15]–[19]. These approaches typically extract hand-crafted features from the training set, and subsequently learn the metrics. From this perspective, metric learning essentially performs feature selection when learning the discriminative models. However, these approaches optimize two key components separately or sequentially. If feature representation is not reliable, some useful information is lost in the first step, and we cannot expect that the learnt metric

obtained in the second step will have desirable performance. Hence, it will be a better choice to jointly learn feature representations and the metrics.

Different from those approaches, we propose a novel deep ranking algorithm for person re-identification. Instead of learning a metric over hand-crafted features, our approach learns joint representations and similarities for image pairs directly from the raw image pixels in a unified framework. Person re-identification can be cast as a retrieval problem: given one or more images of an unknown target, the re-identification task is to rank all individuals from the gallery according to their similarities to that target. We follow the principle to position the correct match of a probe at the top of the list against the gallery set. Hence, we penalize any violation of the ranking order by minimizing the cost corresponding to the sum of the rank of the true match of each probe. Leveraging the close connection between evaluation metrics for learning to rank and loss functions for classification [20], we formulate the person re-identification (ranking) task as a seemingly unrelated binary classification problem.

Inspired by its outstanding performance on numerous traditional computer vision tasks [21]–[28], we utilize the deep convolutional neural network (CNN) to build the relation between a pair of pedestrian images and its similarity score. More specifically, our ranking model is built upon a CNN in which feature representation and metric learning are seamlessly integrated. Rather than employ the Euclidean or cosine distance between the features of a pair of images as the metric, we learn the joint representation of that pair and return a similarity score directly. At the training stage, we organize the labeled data into ranking units, each of which consists of a probe, its true match and corresponding reference set. Our deep network then learns a transformation that tends to assign the highest similarity score to the true match in each ranking unit.

Comprehensive evaluations and comparisons clearly demonstrate the marked superiority of our proposed approach over state-of-the-art person re-identification methods. To the best of our knowledge, there are currently two CNN-based person re-identification algorithms: Deep Metric Learning (DML) [29] and deep Filter Pairing Neural Network (FPNN) [30]. Although our approach is not the first to address the person re-identification problem with deep learning, it is more suitable for re-identification and achieves better performance than either of them.

The contributions of this paper can be summarized as follows.

- It proposes a unified deep ranking framework for person re-identification that directly predicts the similarity of a pair of pedestrian images via joint representation learning. There is no need to explicitly design feature representations, matching models or pre-processing. Our approach is more natural than previous re-identification approaches, including both traditional and deep learning based algorithms.
- An effective learning-to-rank algorithm is presented and integrated with the CNN. It penalizes ranking disorders

in the gallery set, and tends to place the true match at the top.
- Extensive evaluation and analysis of the experimental results demonstrate the effectiveness of our approach. We carefully analyze each component of the framework for a fair self-evaluation, and further discuss key elements that may improve performance in a re-identification framework.

In the next section, we review the related work. We then present our proposed approach in Section III, followed by its optimization in Section IV. Section V presents an extensive comparison with state-of-the-art algorithms, and we analyze each component of our method. Section VI concludes the paper and discusses the future work.

## II. RELATED WORK

We review two streams of related work in terms of the technical components of this work: person re-identification and deep representation learning.

### A. Person Re-identification

Many recent studies have addressed the person re-identification problem. Most of them focus primarily on either new descriptors or metric learning for person re-identification.

The aim of person re-identification descriptors is to generate discriminative signatures for pedestrians. Gray *et al.* [9] defined a feature space consisting of raw color channels in numerous color spaces and texture information captured by Gabor and Schmid filters, ensuring that Ensemble of Localized Features (ELFs) carrying more discriminative information were selected by boosting. Tahir *et al.* [31] proposed a cost-and-performance-effective (CoPE) feature selection approach to identify both well-performing and cost-effective feature subset for person re-identification. Faranzana *et al.* [10] proposed the Symmetry-Driven Accumulation of Local Features (SDALF) that exploited the symmetry property of a person through obtaining head, torso, and legs positions to handle view variations. Cheng *et al.* [32] extended the Pictorial Structure (PS) with their Custom Pictorial Structure (CPS) model to estimate body configurations and extract features from each body part. Ma *et al.* [33] developed the BiCov descriptor based on Gabor filters and the covariance descriptor to handle illumination variations. Kviatkovsky *et al.* [34] developed an invariant intra-distribution structure of color under a wide range of imaging conditions. Yang *et al.* [35] employed color naming and proposed the semantic Salient Color Names based Color Descriptor (SCNCD) that demonstrated robustness to photometric variance. However, descriptors of visual appearance are highly susceptible to cross-view variations due to the inherent visual ambiguities and disparities caused by different view orientations, occlusions, illumination and background clutter. It is difficult to achieve a balance between discriminative power and robustness. In addition, some of these methods rely heavily on foreground segmentations, for instance, [10] needs high-quality silhouette masks for symmetry-based partition, and [34] extracts color intra-distribution signatures only from foreground regions.

Metric learning approaches to re-identification usually follow a similar pipeline: extracting features for each image first, and then learning a metric with which the training data have strong inter-category differences and intra-category similarities. Prosser *et al.* [36] developed an ensemble RankSVM to learn a subspace where the potential true match is given the highest ranking. Mignon *et al.* [4] proposed the Pairwise Constrained Component Analysis (PCCA) to learn a projection into a low dimensional space in which the distance between pairs of data points respects the desired constraints, exhibiting good generalization properties in the presence of high dimensional data. In [15], a metric learning framework is used to obtain a robust Mahalanobis metric for Large Margin Nearest Neighbor classification with Rejection (LMNN-R). Zheng *et al.* [6] proposed the Relative Distance Comparison (RDC) approach to maximize the likelihood of a pair of true matches having a relatively smaller distance than that of a wrongly matched pair in a soft discriminant manner. In [16], the authors introduced the KISSME method from equivalence constraints based on a statistical inference perspective. Li *et al.* [8] partitioned the image spaces of two camera views into different configurations according to the similarity of cross-view transforms, and learned different metrics for different locally aligned common space. In [18], Li *et al.* developed a Locally-Adaptive Decision Function (LADF) that jointly learned the distance metric and a locally adaptive thresholding rule. Pedagadi *et al.* [19] utilized the Local Fisher Discriminant Analysis (LFDA) [37] to learn a subspace to reduce the dimensionality of the extracted high dimensional features. Xiong *et al.* [5] further proposed and evaluated the performance of regularized PCCA (rPCCA), kernel LFDA (kLFDA) [37] and Marginal Fisher Analysis (MFA) [38] with different features and kernels. Other methods that deserve mentioning include salience matching [14] and mid-level filter learning [39].

The methods discussed in this subsection share three main drawbacks: (1) their performance is largely limited by the representation power of the hand-crafted features; (2) feature extraction and metric learning are considered as two independent components and optimized separately, so the interaction between them is not well explored; and (3) the learnt metrics are fitted exclusively to the current scenario (dataset), and cannot be generalized to a new scenario without a significant deterioration in performance.

### B. Deep Learning

Recently, approaches that extract features with deep learning structures, the deep CNNs in particular, have shown great potential in various computer vision tasks, including image classification [22], object detection [25], face verification [26], salient object detection [27], and pose estimation [28]. Although deep learning for re-identification has not been fully investigated, the following works are close to our work in the spirit of learning image similarity or ranking. Hu *et al.* [40] presented a new Discriminative Deep Metric Learning (DDML) method for face verification in the wild, and Wu *et al.* [41] employed deep learning architecture to learn a ranking model for image retrieval. However, they learned

deep networks from hand-crafted features. Wang *et al.* [42] proposed a deep ranking model with multi-scale CNNs to learn fine-grained image similarity directly from the image pixels.

To our knowledge, two deep learning based person re-identification algorithms have been proposed. Yi *et al.* [29] utilized a Siamese CNN with a symmetry structure comprising two sub-nets connected by a cosine layer, and proposed a DML approach for re-identification. Given a pair of images, the deep network extracts features of each image independently, and then uses their cosine distance as the metric. Li *et al.* [30] designed an FPNN that takes two images of pedestrians as input and determines whether they have the same identity. The notable difference between these two algorithms is that the FPNN learns the joint representation of two images, while the DML does not. However, learning a network for binary classification does not seem to be a good choice, because positive pairs are much fewer than negative pairs, and thus the learned network tends to predict most input pairs as negative ones due to the great imbalance of training data [43].

To address these problems, we propose a unified deep learning-to-rank framework that learns joint representation and similarities of image pairs directly from image pixels.

## III. DEEP RANKING FRAMEWORK

### A. Overview

In this section, we describe the proposed approach in detail. Figure 2 gives an illustration of our proposed framework. At the training stage, the labeled data are organized into ranking units and then fed into the deep CNN. The CNN is utilized to model the transformation $f(\cdot, \cdot)$ from a pair of pedestrian images to its similarity score. Since the correct match should be positioned at the top of the gallery, we penalize ranking disorders by minimizing the sum of the ranks of positive pairs in each ranking unit. We formulate these two components into our deep ranking framework and perform joint optimization. The learnt CNN conducts similarity computing in one shot at the test time.

### B. Formulation

Before delving deeper into the formulation, we describe some of the terminologies associated with our problem that will be used later. Without loss of generality, let us consider solving the following person re-identification problem in a single-shot case for convenience. Suppose that we are given a training set $\mathcal{X} = \left\{ \left( x_i^A, x_i^B \right) | i = 1, 2, ..., N \right\}$, where $\left( x_i^A, x_i^B \right)$ is a pair of images of the *i-th* person captured by cameras $A$ and $B$, respectively, and $N$ is the number of pedestrians. For a probe image $x$ to be matched against gallery set $\mathcal{G}$, a ranking list should be generated according to the similarity between $x$ and each image in $\mathcal{G}$. There exists only one correct match $x^+$, which should be placed in the top rank by the learnt ranking model. All other samples in the gallery space are considered to be negative matches, denoted by $\mathcal{G}^-$.

Intuitively, if the learnt ranking model is perfect, the correctly matched pair will be assigned a higher similarity score than a mismatched one, which can be expressed as

$$f\left( x, x^+ \right) > f\left( x, y \right), \forall y \in \mathcal{G}^-, \qquad (1)$$
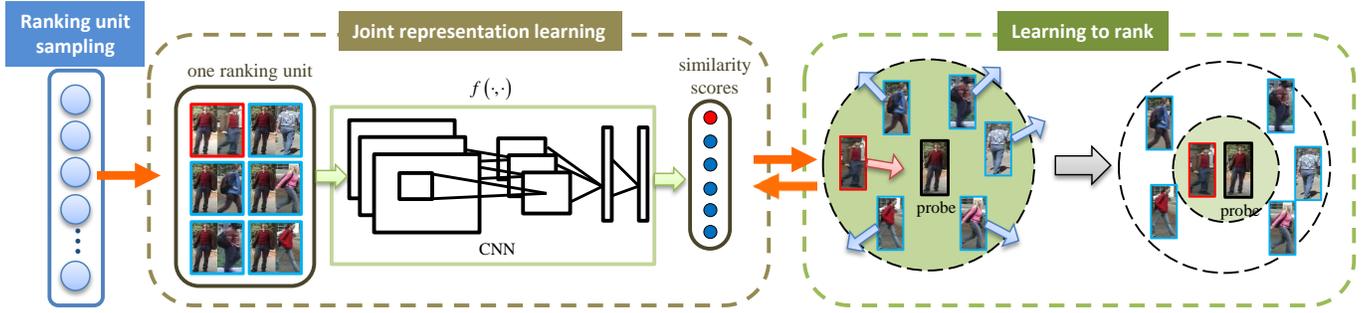
Fig. 2. Illustration of our proposed deep ranking framework, which comprises two key components: deep joint representation learning and learning to rank. We aim to learn a deep CNN that assigns a higher similarity score to the positive pair (marked in red) than any negative pairs (marked in blue) in each ranking unit. Best viewed in color.

where $f(\cdot,\cdot) : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is the learnt similarity metric for an image pair. The rank of $x$ with respect to $\mathcal{G}^-$ can be expressed as a sum of the 0-1 loss function:

$$\text{rank}\left(x|\mathcal{G}^-\right) = \sum_{y\in\mathcal{G}^-} I\left\{f\left(x,x^+\right) - f\left(x,y\right) < 0\right\}, \quad (2)$$

where $I(\cdot)$ is an indicator function whose value is 1 when the expression is true, and 0 otherwise. We propose our learning-to-rank framework based on two main considerations. First, our aim is to place the true match, $x^+$, at the top with regard to $\mathcal{G}^-$. In other words, $\text{rank}\left(x|\mathcal{G}^-\right)$ has to be small. Second, for two mismatches, $y_i, y_j \in \mathcal{G}^-$, we have no idea which is more similar to a given probe $x$, and simply ignore the intra-ranking orders of $\mathcal{G}^-$. Therefore, we formulate the objective function as follows

$$\begin{aligned} J &= \sum_x \text{rank}\left(x|\mathcal{G}^-\right) \\ &= \sum_x \sum_{y\in\mathcal{G}^-} I\left\{f\left(x,x^+\right) - f\left(x,y\right) < 0\right\}. \end{aligned} \quad (3)$$

This formulation minimizes the cost corresponding to the sum of the gallery ranking disorders of each probe.

Unfortunately, dealing directly with the 0-1 loss function leads to a non-differentiable optimization problem. The most common solution to this problem is to upper-bound the 0-1 loss by an easy-to-optimize function. Inspired by [20], we utilize logistic loss function $\sigma(x) = \log_2\left(1 + 2^{-x}\right)$ to replace $I\{x < 0\}$, and rewrite the objective function as

$$J = \sum_x \sum_{y\in\mathcal{G}^-} \sigma\left(f\left(x,x^+\right) - f\left(x,y\right)\right). \quad (4)$$

In this model, the most critical component is learning similarity metric $f(\cdot,\cdot)$. Conventional methods usually design hand-crafted features and subsequently learn a Mahalanobis metric to maximize the inter-class variations and minimize the intra-class variations. In this work, we propose to take advantage of deep CNNs to learn $f(\cdot,\cdot)$ directly from raw image pixels rather than from hand-crafted features. In the following subsection, we introduce the deep network architecture used in our ranking framework.

### C. Network Architecture

Our deep network learns mapping $f(\cdot,\cdot)$ from two images to their similarity score directly. It comprises five convolutional layers to extract features hierarchically, followed by three fully connected layers. Figure 3 shows the detailed structure of our network, which is similar to the popular *AlexNet* [22]. We propose to learn joint representation for an image pair (explained in Section V-F1). Therefore, a notable difference is that we simply stitch two pedestrian images horizontally to form an image that is used as input. More specifically, the images in the pair are both resized to $H \times W$ (here, $H = 2W$, and we set $H = 256$ in the experiments) and then stitched together to form a square image for input. This approach ensures that the aspect ratio of the images remains nearly unchanged, and we do not need to design a new network architecture with two entrances. The convolution operation is expressed as

$$\mathbf{x}_i^{(l)} = \text{relu}\left(b_i^{(l)} + \sum_j \mathbf{k}_{ij}^{(l)} \otimes \mathbf{x}_j^{(l-1)}\right), \quad (5)$$

where $\mathbf{x}_i^{(l)}$ and $\mathbf{x}_j^{(l-1)}$ denote the $i$-th output channel at the $l$-th layer and the $j$-th input channel at the $(l-1)$-th layer, respectively; $\mathbf{k}_{ij}^{(l)}$ is the convolutional kernel between the $i$-th and $j$-th feature map; and $b_i^{(l)}$ is the bias of the $i$-th map. The Rectified Linear Unit (ReLU) is used as the neuron activation function, denoted as $\text{relu}(x) = \max(x, 0)$. Max-pooling is utilized in the first, second and fifth convolutional layers, formulated as

$$\mathbf{x}_{(i,j)}^{(l)} = \max_{\forall(p,q)\in\Omega_{(i,j)}} \mathbf{x}_{(p,q)}^{(l)}, \quad (6)$$

where $\Omega_{(i,j)}$ stands for the pooling region with index $(i,j)$. Since the variations in poses and viewpoints always appear, max-pooling enhances the robustness to small translations [44]. Max-pooling at the first two layers is followed by local response normalization, leading to feature maps that are robust to illumination and contrast variations.

The last three layers are fully connected, expressed as

$$\mathbf{x}^{(l)} = \mathbf{w}^{(l)} \cdot \mathbf{x}^{(l-1)} + \mathbf{b}^{(l)}, \quad (7)$$

where $\mathbf{w}^{(l)}$ and $\mathbf{b}^{(l)}$ are the weight and bias, respectively. The first two fully connected layers reduce the dimensionality of
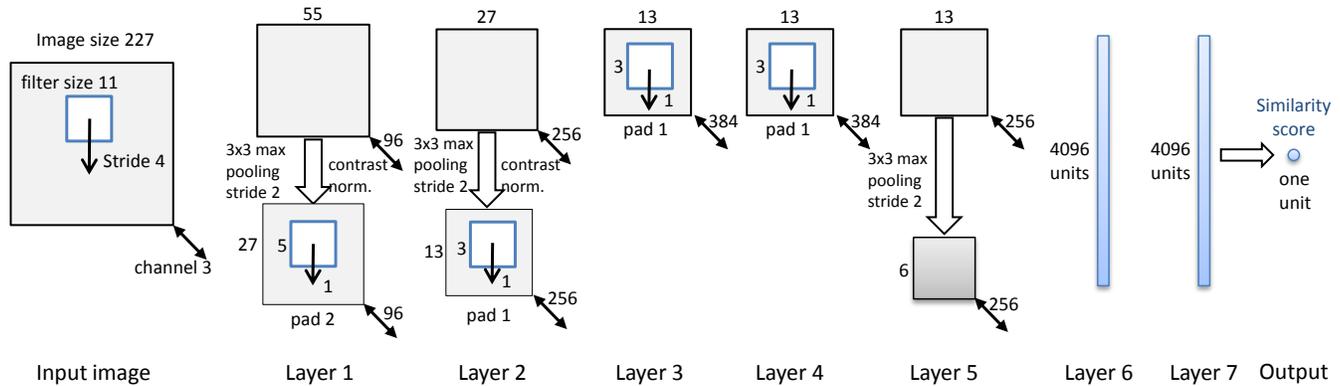
Fig. 3. Architecture of our deep network. A pair of three-channel pedestrian images is first stitched, and then a $227 \times 227$ random crop is presented as the input, which is convolved with 96 different first layer filters, each of size $11 \times 11$, using a stride of 4 in both $x$ and $y$. The resulting feature maps are then passed through a rectified linear unit (ReLU; not shown in this figure), max-pooled ($3 \times 3$ regions with stride 2), and contrast normalized across the feature maps to give 96 different $27 \times 27$ feature maps. Similar operations are repeated in the second to fifth layers. The last three layers are fully connected, taking features from the top convolutional layer as the input in vector form. Finally, a similarity score for the pair is returned.

the extracted joint features from 9216 ($6 \times 6 \times 256$) to 4096, and form highly compact and predictive features, denoted by $\phi(x, y)$. The last layer acts as the similarity/distance metric for $\phi(x, y)$, which can be expressed as

$$f(x, y) = \langle \phi(x, y), \mathbf{w} \rangle + b, \qquad (8)$$

where $f(\cdot, \cdot)$ denotes the similarity metric as before, and $\langle \cdot, \cdot \rangle$ denotes the inner-product between the two vectors. Our network is capable of jointly learning the features and similarity metric with supervised similarity information provided by the proposed ranking algorithm, which characterizes the relative similarity ranking orders.

Note that the activation function for all layers (except the last layer) is the ReLU. Dropout [22] is used in the first two fully connected layers to alleviate over-fitting.

Our reasoning for employing this very deep network architecture for person re-identification is as follows. Since the appearance of a given pedestrian undergoes drastic changes due to complex variations in illumination, pose, viewpoint, camera setting and background clutter across camera views, we argue that the network should be deep enough to handle the inherent visual ambiguities. Further, deeper network learning requires more training samples, but we are given only small-scale labeled data, particularly in a single-shot modality. For instance, the well-known VIPeR dataset [9], which has only 1,264 images for 632 subjects, is far from sufficient for network learning. It seems that the depth of the network is necessarily limited by the amount of training data. In the experiment, we show that this dilemma can be resolved by pre-training and other strategies, as explained in Section V.

## IV. OPTIMIZATION

Our network is trained using the stochastic gradient descent (SGD) algorithm with momentum. Training data $\mathcal{X}$ are organized into mini-batches consisting of several ranking units. The training errors are computed for each mini-batch, and back-propagated to the lower layers.
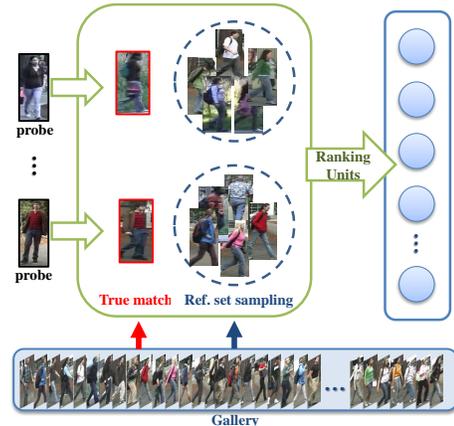


Fig. 4. Illustration of ranking unit sampling.

### A. Ranking Unit Sampling

As discussed in the previous section, we organize the training data into ranking units. Here, we take only a subset of $\mathcal{G}^-$, i.e., $\mathcal{R}_x \subseteq \mathcal{G}^-$, into account. It is considered the reference set for probe image $x$. Each unit comprises a probe $x$, its true match $x^+$, and its corresponding reference set $\mathcal{R}_x$. Note that the reference set $\mathcal{R}_x$ is randomly sampled given probe $x$ (Figure 4). We consider $\mathcal{R}_x$ alone rather than the full $\mathcal{G}^-$ for three reasons: (1) there is no need to load all data into memory at the training stage if we simply sample a subset of $\mathcal{G}^-$, making it more practical for large-scale learning; (2) since reference set $\mathcal{R}_x$ is randomly sampled, the training data of each batch possess a high degree of diversity, which is of great importance to learning; and (3) dealing with a random subset $\mathcal{R}_x$ in each iteration is approximately equivalent to taking the whole of $\mathcal{G}^-$ into account with sufficient iteration times (as explained in Section V-F4). Therefore, the loss is formulated as

$$\mathcal{L} = \sum_x \sum_{y \in \mathcal{R}_x} \sigma \left( f\left(x, x^+\right) - f\left(x, y\right) \right). \qquad (9)$$

The gradients of the ranking loss with respect to the similarities within a ranking unit are

$$\frac{\partial \mathcal{L}}{\partial f(x, x')} = \begin{cases} \dfrac{\delta(x', x^+|x)}{1 + \delta(x', x^+|x)} & x' \in \mathcal{R}_x \\ \displaystyle\sum_{y \in \mathcal{R}_x} \dfrac{-\delta(y, x^+|x)}{1 + \delta(y, x^+|x)} & x' = x^+ \end{cases} \quad (10)$$

where $\delta(i, j|x) = 2^{f(i,x) - f(j,x)}$. The back-propagation algorithm adjusts $f(\cdot, \cdot)$ such that $f(x, x^+)$ is assigned the highest similarity score in the corresponding ranking unit.

Let $\mathcal{X}^A$ and $\mathcal{X}^B$ denote the training data captured by cameras *A* and *B*, respectively. For a probe image from one camera (say, camera *A*), $x_i^A$, the correct match should be $x_i^B$. Here, the reference set can be expressed as $\mathcal{R}_{x_i^A} \subseteq \mathcal{G}_{x_i^A}^-$, where $\mathcal{G}_{x_i^A}^- = \{x_j^B | x_j^B \in \mathcal{X}^B, j \neq i\}$. When training our deep network, we set $|\mathcal{R}_x| = 1$ at the beginning, where $|\cdot|$ is the cardinality of a set. Note that the ranking unit now degrades into a simple triplet constraint, i.e., $f(x, x^+) > f(x, x^-)$ for a triplet $t = (x, x^+, x^-)$. As the learning procedure progresses, we gradually increase the cardinality of the reference set $|R_x|$ in each mini-batch up to 4. From another perspective, the positive pairs and the same number of negative pairs are fed into the deep network at the beginning for balancing. It is not a complicated task to satisfy the triplet constraints tentatively. Since there are much more negative pairs than positive ones, we gradually increase the number of negative samples up to a ratio of 4:1. Now, the problem becomes increasingly difficult because we want the ranking model to position the correct match at the top against a reference set that is increasing in size. In this way, a discriminative ranking model is obtained.

### B. Training Strategies

Several critical training strategies are discussed in this subsection.

**Pre-training** - As previously discussed, a large amount of training data is needed for learning because of the great depth of our network. Since the available labeled data at hand are scarce, we use the labeled data collected from other scenarios (datasets) to assist network learning even though they are subject to quite different distributions. In our experiment, we first use large-scale data to learn a pre-trained model. For each specific scenario (dataset), we initialize the parameters with that pre-trained model, and then fine-tune all layers by back-propagation through the whole network with the given training set. We find experimentally that pre-training is a critical component that boosts performance significantly.

**Relaxing the cross-view constraint** - Recall that the reference set $\mathcal{R}_{x_i^A}$ for probe image $x_i^A$ is rigorously sampled from $\mathcal{G}_{x_i^A}^- = \{x_j^B | x_j^B \in \mathcal{X}^B, j \neq i\}$. In a single-shot modality, the data are insufficient to construct a reference set exhibiting strong diversity. Intuitively, distinguishing two persons from the same camera is relatively easy, and it also helps to learn the similarity metric. Under this consideration we relax the cross-view constraint and sample the reference set from both camera views, i.e., $\mathcal{G}_{x_i^A}^- = \{x_j | x_j \in \mathcal{X}, j \neq i\}$.

**Data augmentation** - It is a common method to artificially enlarge the training data using label-preserving transformations to reduce over-fitting [22], [23]. We also employ similar data augmentation in the form of random crops and horizontal flips. The only notable difference is that the flips are not generated by flipping the images around the $y$-axis directly. Since our input is essentially a pair of images, and we thus flip the image in a different way. Each sub-image is flipped around its own central vertical axis with probability 0.5, and the two sub-images further exchange their positions, also with probability 0.5, which increases the size of our training set by a factor of 8. We perform random crops by randomly extracting $227 \times 227$ patches from the original images (or their horizontal reflections), and then train our network on the extracted patches. At the test time, we deterministically extracted the central crop of the image in addition to its horizontal reflections, and returned the average of these eight scores. We also tried to extract five patches, including the center and four corner patches, as in [22], but achieved similar results.

## V. EXPERIMENT

In this section, we report the results of extensive experiments carried out to compare our approach with state-of-the-art approaches, including both traditional and deep learning based methods, and to evaluate each component of our method in detail. Although the superiority of our approach comes from the framework as a whole, we also carefully assessed each component to give a fair self-evaluation.

### A. Experimental Settings

**Datasets** - To validate our approach, we performed experiments on three benchmark datasets: the VIPeR dataset [9], the CUHK-01 dataset [17], and the CAVIAR4REID dataset [32]. These datasets are highly challenging because of the different camera settings, geometric deformations and photometric variations in different views (Figure 5).

**Evaluation protocol** - We adopted the single-shot modality on the VIPeR and CUHK-01 datasets, and multi-shot modality (with both $N = 5$ and $N = 10$) on the CAVIAR4REID dataset as most previous studies did to allow extensive comparison. Following the commonly used evaluation protocol in [9], we randomly partitioned the dataset into two parts, one half for training and the other for testing, without any overlap in person identities. Each probe image was matched against the gallery set, and the rank of the true match was obtained. The rank-*k* recognition rate is the expectation of the matches at rank *k*, and the cumulative values of the recognition rate at all ranks were recorded as the one-trial Cumulative Matching Characteristic (CMC) result. We repeated the evaluation ten times, and here report the average CMC curve to achieve stable statistics. Note that the CUHK-01 dataset has more than one image of each person, and thus we randomly select one to form the gallery. For the CAVIAR4REID dataset, we used all ten images per view for multi-shot setting with $N = 10$, and randomly select five images for the setting with $N = 5$.

(a) Samples from the VIPeR dataset [9]



(b) Samples from the CUHK-01 dataset [17]



(c) Samples from the CAVIAR4REID dataset [32]

Fig. 5. Samples of pedestrian images observed in different camera views. Each column corresponds to the same identity.

**Implementation details** - In our experiment, the CUHK-02 dataset [8] was used to learn a pre-trained model. Note that the CUHK-02 dataset contains five pairs of views (P1-P5), and P1 is also called the CUHK-01 dataset. The samples from P1 (i.e., the CUHK-01 dataset) were excluded when learning the pre-trained model, because the CUHK-01 was used for evaluation. This ensures that no sample from the test set was used during the pre-training stage by mistake. We implemented our model under the open source Caffe CNN library [45], and trained it using the SGD with momentum of 0.9, weight decay of 0.0005, and learning rates of $10^{-4}$. The parameters were initialized with the model in [22]. We fixed the mini-batch size of 16 ranking units. The positive-negative ratio is set as 1:1 initially. We gradually increased the size of reference set $\mathcal{R}_x$ to 2 and 4, whilst the ratio changes to 1:2 and 1:4 accordingly (see Figure 4).

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

In this section, we compare our proposed method with the following state-of-the-art approaches: ELF [9], SDALF [10], LMNN [46], ITML [47], eSDC [13], Generic Metric [17], Salience Matching (SalMatch) [14], Mid-Level Filter (MLF) [39], eBiCov [33], PCCA [4], LF [19], LADF [18], MFA [5], kLFDA [5], rPCCA [5], RDC [6], attribute-based PRDC (aPRDC) [11], CPS [32], RankSVM [36], LMNN-R [15], KISSME [16], SCNCD [35], ICT [48], and Feature Warps (FW) [49]. Note that not all of these approaches have reported

results for all three datasets. For instance, VIPeR is the most widely used benchmark, and thus several researchers have reported the results of various approaches on VIPeR but not CUHK-01 or CAVIAR4REID. Here, we compare our method with the foregoing methods if available.

*1) Performance on the VIPeR dataset:* The VIPeR dataset [9] contains 632 person image pairs. Each pair has two images of the same person observed from different views, resized to $128 \times 48$. Most of the approaches considered have CMC curves reported for this dataset, hence we give a more detailed comparison for VIPeR. Figure 6(a) shows the CMC curves up to rank 25 comparing our method with state-of-the-art methods. It is obvious that our method gives the best result in the main. To present the quantized comparison results more clearly, we also summarize the performance comparison at several top ranks in Table I. It can be seen that, our proposed method achieves a 38.37% rank-1 matching rate, outperforming the previous best result that of SCNCD [35], which achieved 37.8%. The other best performing methods on the VIPeR dataset are metric learning algorithms such as kLFDA and MFA [5]. Our method performs best over ranks 1, 5, and 10, whereas kLFDA is the best at rank-20. Our experimental results suggest that even though our model suffers from a severe lack of training data, it still achieves state-of-the-art performance on the highly challenging VIPeR dataset.

Note that several of the methods considered have been combined with other descriptors to achieve better performance, such as eSDC [13] and eBiCov [33]. To the best of our knowledge, the current best result on the VIPeR dataset is that achieved by a combination of MLF and LADF [39]. Leveraging the considerable complementarity between the traditional framework and deep networks [50], we also report the results of our approach in combination with existing metric learning approaches with hand-crafted low-level features. Here, we simply sum up the scores of our method and kLFDA under the same training/testing partitions, and recompute the CMC curve. As shown in Figure 6(a) and Table I, the rank-1 matching rate surges to about 53%, far surpassing all of the state-of-the-art methods considered.

*2) Performance on the CUHK-01 dataset:* The CUHK-01 dataset [17] is larger in scale than the VIPeR. It contains 971 persons, each of whom has two images in each camera view. Camera *A* captures the frontal or back views of pedestrians, whereas camera *B* captures their side views. All images are normalized to $160 \times 60$. Note that we used two images of each person for training, and randomly selected only one for the test.

We compared our proposed method with several state-of-the-art approaches, such as MLF [39] and SalMatch [14]. The CMC curves obtained using the $\ell_1$-norm and $\ell_2$-norm distances of concatenated dense features were also compared as baselines [14]. As shown in Figure 6(b) and Table II, our method outdistances all state-of-the-art methods at all ranks, which again validates its effectiveness. Our method achieves a rank-1 matching rate of 50.41%, outperforming the previous best result reported by MLF, which achieved a 34.30% rank-1 matching rate, by a sizeable margin. The significant advantage
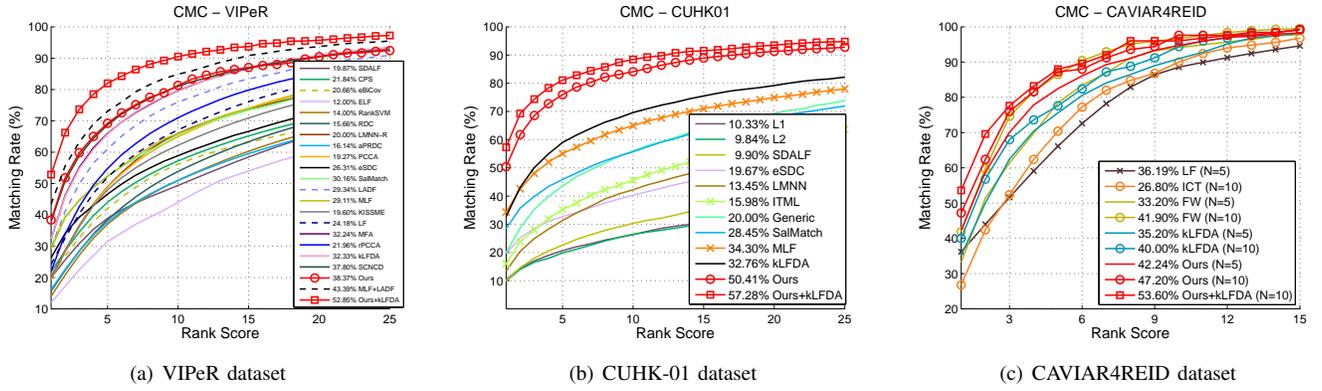
Fig. 6. Performance comparisons with state-of-the-art approaches using CMC curves on the VIPeR ($p = 316$), CUHK-01 ($p = 486$) and CAVIAR4REID ($p = 25$) datasets. In the legends, we also report the rank-1 matching rate for each approach. Best viewed in color.

| Method | $r = 1$ | $r = 5$ | $r = 10$ | $r = 20$ |
|---|---|---|---|---|
| ELF [9] | 12.00 | 41.50 | 59.50 | 74.50 |
| SDALF [10] | 19.87 | 38.89 | 49.37 | 65.73 |
| CPS [32] | 21.84 | 44.00 | 57.21 | 71.00 |
| RDC [6] | 15.66 | 38.42 | 53.86 | 70.09 |
| aPRDC [11] | 16.14 | 37.72 | 50.98 | 65.95 |
| RankSVM [36] | 14.00 | 37.00 | 51.00 | 67.00 |
| KISSME [16] | 19.60 | 48.00 | 62.20 | 77.00 |
| PCCA [4] | 19.27 | 48.89 | 64.91 | 80.28 |
| rPCCA [5] | 21.96 | 54.78 | 70.97 | 85.29 |
| eBiCov [33] | 20.66 | 42.00 | 56.18 | 68.00 |
| LMNN-R [15] | 20.00 | 49.00 | 66.00 | 79.00 |
| eSDC [13] | 26.31 | 46.61 | 58.86 | 72.77 |
| SalMatch [14] | 30.16 | 52.31 | 65.54 | 79.15 |
| MLF [39] | 29.11 | 52.34 | 65.95 | 79.87 |
| LF [19] | 24.18 | 52.00 | 67.12 | 82.00 |
| LADF [18] | 29.34 | 61.04 | 75.98 | 88.10 |
| MFA [5] | 32.24 | 65.99 | 79.66 | 90.64 |
| kLFDA [5] | 32.33 | 65.78 | 79.72 | **90.95** |
| SCNCD [35] | 37.80 | 68.67 | 81.01 | 90.51 |
| **Ours** | **38.37** | **69.22** | **81.33** | 90.43 |
| MLF + LADF [39] | 43.39 | 73.04 | 84.87 | 93.70 |
| **Ours + kLFDA** | **52.85** | **81.96** | **90.51** | **95.73** |

TABLE I
TOP-RANKED MATCHING RATES (%) ON THE VIPeR DATASET ($p = 316$).
THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

| Method | $r = 1$ | $r = 5$ | $r = 10$ | $r = 20$ |
|---|---|---|---|---|
| $\ell_1$-norm [14] | 10.33 | 20.64 | 26.34 | 33.52 |
| $\ell_2$-norm [14] | 9.84 | 19.84 | 26.42 | 33.13 |
| SDALF [10] | 9.90 | 22.57 | 30.33 | 41.03 |
| eSDC [13] | 19.67 | 32.72 | 40.29 | 50.58 |
| LMNN [46] | 13.45 | 31.33 | 42.25 | 54.11 |
| ITML [47] | 15.98 | 35.22 | 45.60 | 59.81 |
| Generic Metric [17] | 20.00 | 43.58 | 56.04 | 69.27 |
| SalMatch [14] | 28.45 | 45.85 | 55.67 | 67.95 |
| MLF [39] | 34.30 | 55.06 | 64.96 | 74.94 |
| kLFDA | 32.76 | 59.01 | 69.63 | 79.18 |
| **Ours** | **50.41** | **75.93** | **84.07** | **91.32** |
| **Ours+kLFDA** | **57.28** | **81.07** | **88.44** | **93.46** |

TABLE II
TOP-RANKED MATCHING RATES (%) ON THE CUHK-01 DATASET
($p = 486$). THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

| Method | $r = 1$ | $r = 5$ | $r = 10$ | $r = 20$ |
|---|---|---|---|---|
| FW ($N = 5$) [49] | 33.20 | 78.50 | 94.10 | **100.00** |
| LF ($N = 5$) [19] | 36.19 | 66.15 | 88.56 | 98.41 |
| kLFDA ($N = 5$) | 35.20 | 75.60 | 90.96 | 99.76 |
| **Ours** ($N = 5$) | **42.24** | **82.48** | **94.72** | 99.92 |
| ICT ($N = 10$) [48] | 26.80 | 70.40 | 90.00 | 99.60 |
| FW ($N = 10$) [49] | 41.90 | 86.50 | 96.70 | 100.00 |
| kLFDA ($N = 10$) | 40.00 | 77.60 | 94.40 | 100.00 |
| **Ours** ($N = 10$) | **47.20** | **87.20** | **97.60** | **100.00** |
| **Ours+kLFDA** ($N = 10$) | **53.60** | **88.00** | 96.00 | **100.00** |

TABLE III
TOP-RANKED MATCHING RATES (%) ON THE CAVIAR4REID DATASET
($p = 25$). THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

Note that no previous work has reported the CMC curve of kLFDA on the CUHK-01 dataset, and hence we conducted this experiment with the codes offered by [37] and commonly used low-level features. Similar to [6], we extracted two types of low-level features: color and texture for each image. More specifically, we equally partitioned each image into $N$ horizontal stripes. For each stripe, color hisograms and Gabor texture features were extracted. Each feature channel was represented as an $\ell_1$-normalized 16-bin histogram, and all histograms were concatenated to form a single feature vector. As shown in Figure 6(b) and Table II, the combination of our approach with kLFDA improves the rank-1 matching rate by about 7%.

*3) Performance on the CAVIAR4REID dataset:* The CAVIAR4REID dataset is composed of 1,220 images of 72 pedestrians out of which 50 are viewed by two disjoint cameras. It has broad changes in resolution, the minimum and maximum size of the images is $17 \times 39$ and $72 \times 144$, respectively. It is noteworthy that both of the VIPeR and CUHK-01 are kind of regular, in the sense that pedestrians are rigidly enclosed under fixed-size bounding boxes, while CAVIAR4REID teems with significant variations in image resolution. Thus, we conducted experiment on this dataset in order to verify whether resizing would impact the proposed method.

Following [49], we considered only 50 persons viewed by two cameras, and discarded the remaining 22 persons who appear in only one camera. In this way, the 50 persons are

of our proposed method is that more training data are fed to the deep network to learn a data-driven solution for the specific scenario in question.

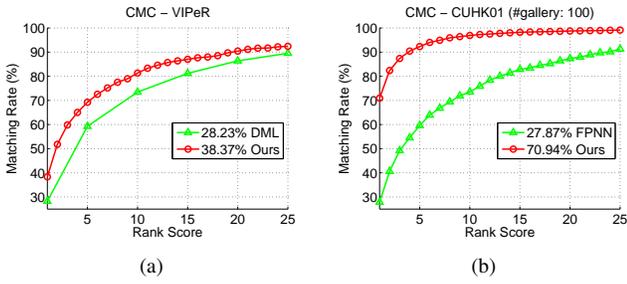As before, we also combined our method with kLFDA.

Fig. 7. Comparison with two deep learning based methods: DML [29] and FPNN [30]. (a) Comparison with DML on VIPeR ($p = 316$); (b) Comparison with FPNN on CUHK-01 ($p = 100$). Note that the gallery contains only 100 subjects, which is different from that in Figure 6(b) and Table II.

equally divided into training and test sets of 25 persons each. We compare our method with state-of-the-art approaches in the multi-shot modality with both $N = 5$ and $N = 10$. Since kLFDA performs better on both VIPeR and CUHK-01 datasets, we conducted additional experiment and also compare the CMC curves obtained by kLFDA as before. Figure 6(c) and Table III present the results, showing that our proposed method outperforms all previous methods in both settings with $N = 5$ and $N = 10$, respectively. Therefore, we can conclude that our proposed method is also robust to severe variation in image resolution. In addition, as expected, the combination of our method with kLFDA boosts the rank-1 matching rate by over 6%.

These comparisons clearly suggest that our proposed method outperforms all state-of-the-art algorithms, particularly when sufficient training data are provided. The main reason for its superior performance is that our framework is capable of jointly tackling representation learning and the learning-to-rank task rather than requiring two-step separate optimization.

## C. Comparision with CNN-based Person Re-identification Algorithms

In this section, we present the results of comparison between our method and two deep learning based person re-identification algorithms: DML [29] and deep FPNN [30]. DML has been used only in experiments on the VIPeR dataset. Li *et al.* [30] concluded that existing datasets are too small to train deep networks, and thus they only conducted their experiments on both the large-scale CUHK-03 dataset and the CUHK-01 dataset. No previous CNN-based algorithm reported the results on the CAVIAR4REID. Therefore, we compare our method with DML on VIPeR and with FPNN on CUHK-01. Table IV(a) and Figure 7(a) suggest that our model significantly surpasses DML, particularly at rank-1 (over 10%). Note that the FPNN experiment on CUHK-01 was conducted in a different setting, with only 100 persons chosen for testing and the remaining 871 persons used for training and validation. Recall that we used 486 persons for testing and only 485 for training. In other words, we used an approximately 5-fold larger gallery than that in [30]. Even though our setting was much more challenging than the FPNN setting, we still achieve a 50.41% rank-1 matching rate, far surpassing that of FPNN, which was only 27.87%. For a

fairer comparison, we randomly removed 386 persons from the gallery and recomputed the average CMC curve, as shown in Figure 7(b) and Table IV(b). Our proposed method achieves a 70.94% rank-1 matching rate, a greater than 43% improvement over FPNN.

Note that DML, FPNN, and our proposed method learn feature representation with deep CNNs, and hence their performance relies primarily on the ranking algorithm rather than the representation power of CNNs. These experimental results show that our ranking mechanism model is more suitable than others for person re-identification.

(a) Comparison with DML on the VIPeR dataset

| Method | $r = 1$ | $r = 5$ | $r = 10$ | $r = 20$ | $r = 30$ |
|---|---|---|---|---|---|
| DML [29] | 28.23 | 59.27 | 73.45 | 86.39 | 92.28 |
| **Ours** | **38.37** | **69.22** | **81.33** | **90.43** | **94.15** |

(b) Comparison with FPNN on the CUHK-01 dataset

| Method | $r = 1$ | $r = 5$ | $r = 10$ | $r = 20$ | $r = 30$ |
|---|---|---|---|---|---|
| FPNN [30] | 27.87 | 59.64 | 73.53 | 87.34 | 93.92 |
| **Ours** | **70.94** | **92.30** | **96.90** | **98.74** | **99.34** |

TABLE IV
COMPARISON WITH TWO DEEP LEARNING BASED PERSON
RE-IDENTIFICATION ALGORITHMS: DML [29] AND FPNN [30]. THE BEST
RESULTS ARE HIGHLIGHTED IN **BOLD**.

## D. Evaluation of Open-World Scenarios

The CMC-based evaluation protocol used in Section V-B and V-C assumes that the gallery and probe sets contain exactly the same individuals. Thus, this can be regarded as the *close-world* re-identification task. Here we also consider the person re-identification problem in the context of open-world scenarios [51]. The main difficulty of open-world setting lies in the fact that the probe image may not belong to anyone on the gallery/target set. Under this setting, a watch-list of a handful of known people is provided as the target set. Additionally, there are a large amount of non-target imposters captured along with the target set. Given a probe image, we need to determine whether it is on the watch-list or not.

Following the open-world evaluation metrics defined in [51], we exploited the True Target Rate (TTR) and False Target Rate (FTR) as,

$$TTR = \frac{\#\mathcal{TTQ}}{\#\mathcal{TQ}}, \quad (11)$$

$$FTR = \frac{\#\mathcal{FNTQ}}{\#\mathcal{NTQ}}, \quad (12)$$

where $\mathcal{TQ}$ denotes the query target images from target people; $\mathcal{NTQ}$ indicates the query non-target images from non-target people; $\mathcal{TTQ}$ denotes the query target images that are verified as one of the target people; and $\mathcal{FNTQ}$ means the query non-target images that are mistakenly verified as one of the target people. In the experiment, we randomly selected $p$ subjects from the gallery set as the target set, and removed the remaining images from the gallery. In this way, most probe images cannot find their true matches in the target set, which contains several pedestrians that we are interested in. As in the close-world setting, we computed the similarity
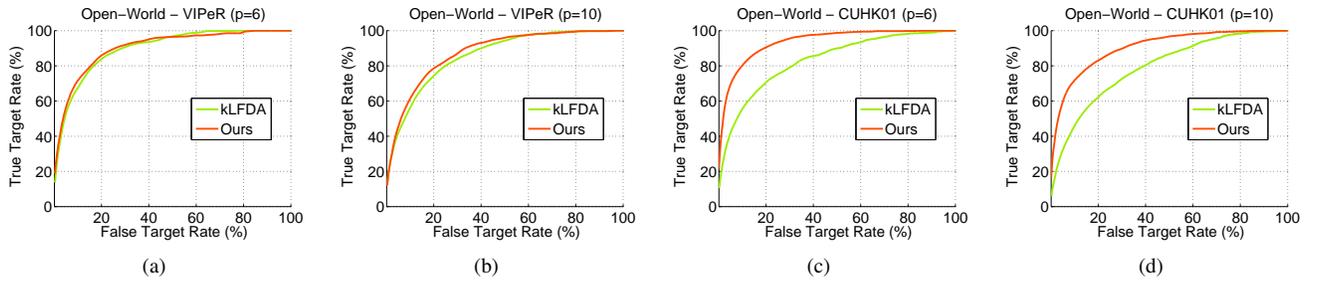
Fig. 8. Open-world set-based verification performance comparisons with kLFDA using TTR-FTR curves. Best viewed in color.

| Dataset | VIPeR ($p = 6$) | | | | | | VIPeR ($p = 10$) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FTR | 0.1% | 1% | 5% | 10% | 20% | 50% | 0.1% | 1% | 5% | 10% | 20% | 50% |
| kLFDA | 14.58 | 23.96 | 52.71 | 67.71 | 83.96 | **96.67** | 12.88 | 18.88 | 42.63 | 56.38 | 74.50 | 94.38 |
| **Ours** | **18.75** | **29.79** | **58.33** | **71.46** | **86.25** | 96.46 | **13.88** | **19.38** | **46.25** | **61.00** | **78.50** | **95.88** |
| Dataset | CUHK-01 ($p = 6$) | | | | | | CUHK-01 ($p = 10$) | | | | | |
| FTR | 0.1% | 1% | 5% | 10% | 20% | 50% | 0.1% | 1% | 5% | 10% | 20% | 50% |
| kLFDA | 10.53 | 19.13 | 42.40 | 55.60 | 70.87 | 89.93 | 6.80 | 12.52 | 32.72 | 46.04 | 62.28 | 86.68 |
| **Ours** | **22.87** | **40.13** | **69.53** | **80.07** | **90.60** | **98.73** | **14.92** | **32.84** | **59.48** | **71.96** | **83.04** | **96.72** |

TABLE V

OPEN-WORLD SET-BASED VERIFICATION RESULTS: TRUE TARGET RATE (TTR) IN % AGAINST FALSE TARGET RATE (FTR). THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

score for each probe image and each target image. A value $s$ is used to threshold these scores and therefore a curve depicting the TTR value against the FTR value is reported by changing the value $s$. We also reported the TTR value when the FTR value is fixed. We repeated this experiment 100 times and reported the average verification performance on the VIPeR and CUHK-01 datasets, respectively. The comparison of our method and kLFDA is given. Figure 8 shows the TTR-FTR curves with $p = 6$ and $p = 10$ (set as that in [51]). More detailed comparison of TTR values with FTR fixed is presented in Table V. It demonstrates that, our proposed method also surpasses kLFDA under the open-world scenarios, particularly when sufficient training samples are given.

### E. Comparison of Performance across Datasets

In this section, we compare the performance in the cross dataset setting which better coincides with practical applications. In the case of large-scale camera networks, it is impossible to collect enough labeled data to learn a specific discriminative model for each camera pair. Therefore, transferring the learnt model to the current scenario without significant degradation in performance is a more practical approach. The Domain Transfer Ranked Support Vector Machine (DTRSVM) [52] was proposed to address this problem. However, the DTRSVM still has to re-train the metric with source domain data and easily collected negative samples from target domain in a multi-task learning framework.

It is known that, deep CNNs have a strong generalization ability across datasets. Here, we compare the cross-dataset performance of our method with that of the DML [29] and the DTRSVM. Following [29], we directly tested their performance on the VIPeR dataset with the pre-trained model learnt from CUHK-02 without fine-tuning. Note that the DTRSVM needs negative samples from VIPeR when
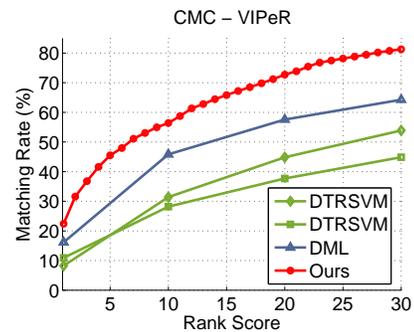


Fig. 9. Comparison of cross-dataset performance. Note that, different from DTRSVM, the DML and our method do not use any sample from VIPeR. Best viewed in color.

transferring the metric, whereas our method and the DML did not use any sample on the VIPeR dataset to fine-tune the network parameters. Figure 9 and Table VI summarize the CMC curves and detailed numerical matching rates, respectively. The experimental results show that: 1) deep CNNs exhibit good generalization ability for cross-dataset re-identification tasks, even though no sample from the target domain is used to fine-tune the networks; 2) using the same training data, our method also outperforms DML in generalization ability significantly, which again demonstrating the superiority of our proposed method.

### F. Evaluations and Analysis

We now analyze a fair self-evaluation of the proposed person re-identification algorithm. Fair self-evaluation is defined as evaluation of both the final output and each component of the algorithm to assess the actual contributions of various components. We are here inspired by [5]. Unlike many works simply comparing the final CMC curves in different experimental

| Method | Training data | need to re-train | $r = 1$ | $r = 10$ | $r = 20$ | $r = 30$ |
|---|---|---|---|---|---|---|
| DTRSVM [52] | i-LIDS + VIPeR | $\sqrt{}$ | 8.26 | 31.39 | 44.83 | 53.88 |
| DTRSVM [52] | PRID + VIPeR | $\sqrt{}$ | 10.90 | 28.20 | 37.69 | 44.87 |
| DML [29] | CUHK | $\times$ | 16.17 | 45.82 | 57.56 | 64.24 |
| **Ours** | CUHK | $\times$ | **22.41** | **56.39** | **72.72** | **81.27** |

TABLE VI
CROSS-DATASET EXPERIMENT: TOP-RANKED MATCHING RATES (%) ON THE VIPeR DATASET. THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

settings, [5] employed uniform and consistent representation in all comparisons to fairly evaluate different metric learning methods, i.e., the important components of re-identification algorithms. It is important to follow this principle, because representation and the ranking mechanism can both affect the final performance. We cannot determine the actual contribution of different ranking mechanisms, only according to the final CMC curves under different representations. Overall, the output result of a re-identification approach is determined by several key factors, mainly including pre-processing, image segmentation, pedestrian representation, ranking mechanism, etc.

In a more general setting than that in [5], a fair self-evaluation for a re-identification algorithm should verify that the effectiveness of the proposed algorithm stems primarily from the components that are claimed to be effective, rather than comparing only the final output or a specific component in different settings. Our aim in conducting a fair self-evaluation is to assess each component of our proposed method to prove the positive roles that all components played in our re-identification framework. Following our analysis of the extensive comparison experiments in Section V-B and V-C, which demonstrate the superiority of our approach as a whole, we here evaluate and analyze each component of our deep ranking framework in detail.

*1) Contribution of joint representation learning:* Many previous works, including shallow and deep learning algorithms, share a similar framework: they extract features for two images separately, and then use the Euclidean or cosine distance as the metric. The only difference between them is whether the features are hand-crafted or learned by CNNs. Our approach learns joint representation for two pedestrian images and directly predicts their similarity from raw pixels, which is similar to FPNN in terms of joint representation learning, but differs from it in learning to rank (further explained below). Our approach is motivated by human assessment: when a person assesses whether two images belong to the same pedestrian, he or she puts the two images together and compares clothes and accessories of those depicted. Actually, the discriminative information comes from different parts of the images. For instance, given three pedestrian images, $a$, $b$, and $c$, that look quite similar, we are able to distinguish $a$ from $b$ by the subjects' knapsacks, which have different colors; and $a$ from $c$ by their different shoes. Given probe $a$, the discriminative region comes from the knapsacks compared with $b$, and from the shoes compared with $c$. However, valuable information is often hidden or ignored when features are extracted independently. We argue that a decision is made jointly from two images rather than from two separately generated items of

discriminative information. Therefore, we propose to predict the similarity of two images via joint representation learning with a learned deep network.

As kLFDA achieves the best performance among metric learning algorithms with hand-crafted features for re-identification at present, we compare the ranking examples of our method and those of kLFDA, as shown in Figure 10. The images in the first three rows are taken from the VIPeR dataset, and those in the last three rows come from the CUHK-01 dataset. This comparison also reveals multiple valuable facts about joint representation learning, as follows. (1) Low-level features can easily produce counterintuitive results, whereas joint representation learning captures semantic colors very well. In the first row, given the probe dressing in pink, our method correctly places the true match at rank-1 and also places other persons wearing pink clothes in the top ranks. However, kLFDA mistakenly positions two persons who wear different colors of clothes above the true match. (2) Joint representation learning exhibits superior performance particularly when illumination varies greatly (e.g., row 2). In rows 4 and 5, joint representation learning is able to capture discriminative information from the green or red bag in the presence of both illumination and pose variations, whereas low-level features fail. (3) Low-level features sometimes outperform joint representation learning provided that the degree of cross-view variation is small. For instance, kLFDA performs well in row 3, because the probe shares similar pose and discriminative background to the corresponding gallery image, which just right provide a useful context cue for ranking. (4) Row 6 is extremely challenging because several candidates look nearly the same. Our method still correctly matches the probe at rank-1, suggesting that joint representation learning is capable of mining subtle discriminative information for re-identification.

*2) Ranking versus direct binary classification:* We also performed experiments to evaluate the contribution of our learning-to-rank algorithm. The results reveal the strength of our ranking mechanism relative to other CNN-based methods. We first removed our proposed ranking model, and then employed a softmax layer to replace the last fully connected layer in Figure 3, with the other layers left unchanged. In this way, the deep network was used to assess whether two input images belonged to the same person. In other words, we performed direct binary classification instead of learning to rank for the person re-identification task like FPNN [30]. The experiments were conducted on the CUHK-01 dataset with a positive-negative ratio of 1:1. No pre-training was used. The CMC curves in Figure 11 show that learning to rank consistently surpasses direct binary classification, thereby demonstrating

Fig. 10. Comparison of the ranking examples of kLFDA and our method. In each row, the left-most image is the probe, and the others are the top 15 matched gallery images of kLFDA and our method. The correct match of the probe is highlighted with a red bounding box. Best viewed in color.
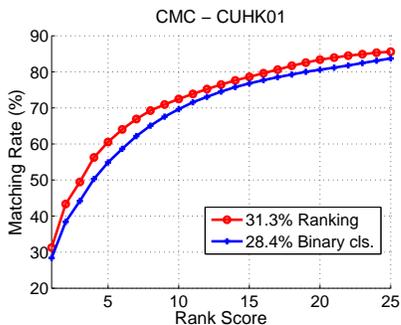


Fig. 11. Comparison of ranking and direct binary classification.

that the good performance of our method stems from both deep representation learning and the ranking algorithm because of the intrinsic difference between image classification and ranking tasks. A good network for image classification is not optimal for ranking [42]. Our learning-to-rank algorithm is based on a relative similarity comparison rather than an absolute decision for each pair, which better accords with the spirit of person re-identification.

*3) Contribution of pre-training:* We evaluated the contribution of pre-training by comparing the performance with and without it. Even though it is difficult to quantify the amount of data needed with reference to the networks depth, deep networks have shown ravenous appetite for training data. As an example, the verification rate improves about 10% for the same deep network, when the number of training images increases from 2.6 million to 26 million [53]. Here, we performed experiments on the CUHK-01 dataset with a positive-negative ratio of 1:1. Figure 12 shows the CMC curves and the loss of the training and test sets. As depicted in Figure 12(a), pre-training improves performance by 15% at rank-1, and

consistently boosts it by about 10% at all ranks. Our deep network needs more training data to learn the parameters, and pre-training is thus able to make use of large-scale outside data. In addition, Figures 12(b) and 12(c) clearly show that the deep network is given better initialization and converges much faster through pre-training. The implication is that the re-identification performance of our proposed method would be further improved through a pre-trained model learned with larger-scale labeled data.

It is noteworthy that pre-training is not an indispensable component to our proposed framework. As discussed in previous section, pre-training allows us to enlarge the depth of our neural network, helping to learn more robust and discriminative representation. In this way, the size of network is no longer limited by the labeled training data collected from the target scenario. As an example, we have to design a small-sized network for the VIPeR dataset without pre-training. Fortunately, leveraging the generalization ability of CNNs, we can borrow outside data to pre-train a deeper network, and then fine-tune its parameters for the current target scenario for better re-identification performance.

However, one issue remains: it seems that the pre-training stage gives the proposed method an "unfair" advantage compared to other state-of-the-art approaches, which have been devised under a constraint of only using the dataset subset for training. To fairly verify the effectiveness of our method, we removed the pre-training stage, trained the network with only the training subset of CUHK-01, and then compared it with FPNN on the CUHK-01 dataset, as shown in Figure 13. The rank-1 matching rate declines from 70.94% to 55.11% when the pre-training stage is excluded. Even so, our method still achieves a greater than 27% improvement over FPNN. It is clear that the remarkable performance of our model
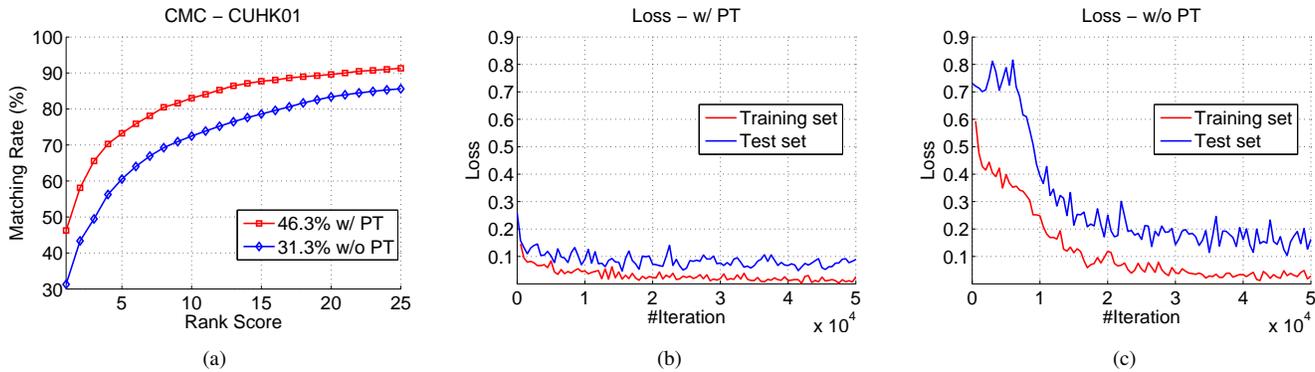
Fig. 12. Analysis of the contribution of pre-training (abbreviated as PT in the figure). These experiments were conducted on the CUHK-01 dataset with a 1:1 positive-negative ratio with or without pre-training. (a) shows the CMC curves, and (b) and (c) show the loss of the training and test sets with and without pre-training, respectively. Best viewed in color.

does not simply come from the advantage of pre-training, although pre-training contributes to better performance indeed. For traditional metric-learning-based algorithms, [52] has experimentally showed that learning metrics from outside data under different distributions leads to severe deterioration in performance. Here we also compare the experimental results when learning kLFDA metrics with different training data to gain insights into how traditional methods are influenced by outside training data. Figure 14 shows the comparisons on the VIPeR and CUHK-01 datasets. Let $\mathcal{S}_1$ denote the original subset for training, and $\mathcal{S}_2$ denotes the data of P2-P5 from the CUHK-02 dataset. In other words, traditional metric learning algorithms learn the models with $\mathcal{S}_1$ that comes from the same dataset. $\mathcal{S}_2$ is exactly the training set that we used to pre-train the deep network. We evaluated three types of training sets: $\mathcal{S}_1$, $\mathcal{S}_2$, and $\mathcal{S}_1 \cup \mathcal{S}_2$. It is demonstrated that (1) learning the metrics from $\mathcal{S}_2$ leads to the worst performance, because there exist significant differences between the distributions of $\mathcal{S}_2$ and the test target sets; (2) little performance gains are observed if we augmented the training set with outside data (i.e., using $\mathcal{S}_1 \cup \mathcal{S}_2$), compared to that learned with $\mathcal{S}_1$; (3) different from traditional metric learning approaches, our proposed method can better leverage the outside data to learn more discriminative representation that possesses strong generalization ability. That is the reason why our proposed method exploits the pre-training strategy.
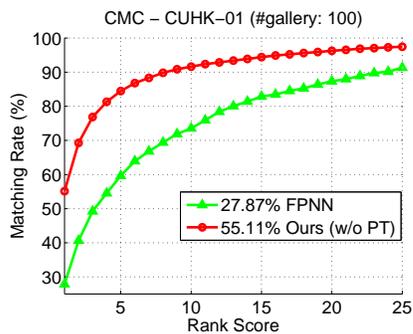


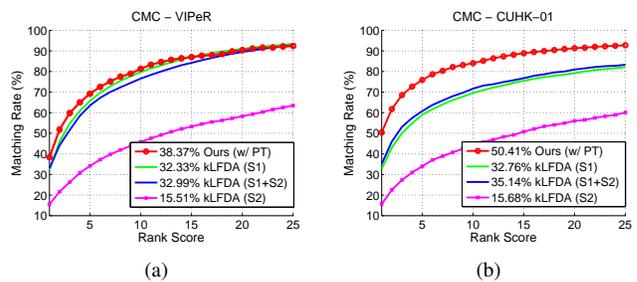Fig. 13. Comparison with FPNN on the CUHK-01 dataset ($p = 100$). Note that no pre-training is used here.



Fig. 14. Comparisons of kLFDA performance with different training data on the VIPeR and CUHK-01 datasets. $\mathcal{S}_1$ denotes the original subset for training, and $\mathcal{S}_2$ denotes the data of P2-P5 from the CUHK-02 dataset. The CMC curves of our method that pre-trains the deep network with $\mathcal{S}_2$ are also presented for comparison.

*4) Analysis of ranking unit sampling:* We here analyze the effectiveness of the ranking unit sampling method. Recall that we initially trained the deep network with $|\mathcal{R}_x| = 1$, i.e., a positive-negative ratio (abbreviated as "ratio" for convenience) of 1:1, and then gradually increased the $|\mathcal{R}_x|$ up to 4 (a ratio of 1:4). Figure 15 shows the CMC curves with different ratios of the reference sets in the ranking units on the VIPeR and CUHK-01 datasets. The CMC curves with a ratio of 1:2 consistently surpassed those with a ratio of 1:1 on both datasets, but no significant improvement was observed when we increased the ratio of the negative pairs to 4. On the VIPeR dataset, the CMC with a 1:4 ratio achieved a better rank-1 matching rate but performed nearly the same after rank-5. On the CUHK-01 dataset, the CMC with a 1:4 ratio worsened slightly but showed better performance after rank-7. We have also conducted more experiments on the change of ratio. The experimental results suggest that increasing $|\mathcal{R}_x|$ gives a small boost in re-identification performance, but the improvement with an increase in $|\mathcal{R}_x|$ is near saturation until 4. We conclude that the randomly sampled ranking units with a ratio of 1:2 can approximately replace the whole gallery set $\mathcal{G}$ during optimization. All of the results compared with the state-of-the-art algorithms had a ratio of 1:2.
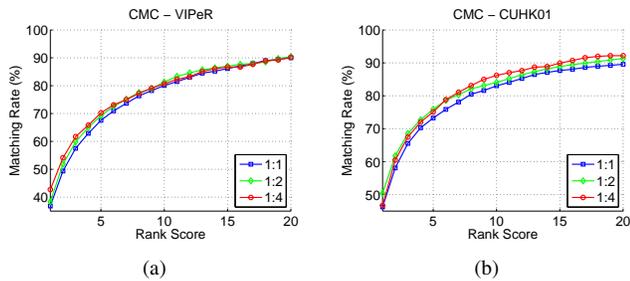
Fig. 15. Analysis of ranking unit sampling. The CMC curves on the VIPeR dataset (a) and CUHK-01 dataset (b) with different positive-negative ratios in the ranking units. Best viewed in color.

## VI. CONCLUSION

In this paper, we formulate the person re-identification task as a learning-to-rank problem, and propose a ranking model that learns a similarity metric that tends to place the true match of a probe at the top by penalizing ranking disorders in the gallery. A deep CNN is utilized to build a relation between image pairs and their similarities. These two components are then seamlessly integrated into a unified deep ranking framework that conducts similarity computing in one shot via joint representation learning directly from raw pixels without feature engineering. Extensive experimental results clearly demonstrate the effectiveness of our proposed approach.

In the future, we would like to further explore ways to make full use of larger-scale outside data for network learning. In addition, we plan to explore how to adapt our approach to video data, i.e., how to measure the similarity of two sequences of detected pedestrian images. Finally, our framework can easily be applied to improve other learning-to-rank tasks such as relative attributes and image interestingness prediction.

## ACKNOWLEDGMENT

## REFERENCES

[1] B. Song, A. T. Kamal, C. Soto, C. Ding, J. A. Farrell, and A. K. Roy-Chowdhury, "Tracking and activity recognition through consensus in distributed camera networks," *IEEE Trans. Image Process.*, vol. 19, no. 10, pp. 2564–2579, Oct. 2010.

[2] R. Vezzani, D. Baltieri, and R. Cucchiara, "People reidentification in surveillance and forensics: A survey," *ACM Comput. Surv.*, vol. 46, no. 2, pp. 29:1–29:37, Nov. 2013.

[3] S. Gong, M. Cristani, C. C. Loy, and T. M. Hospedales, "The re-identification challenge," in *Person Re-Identification*. Springer, 2014, pp. 1–20.

[4] A. Mignon and F. Jurie, "PCCA: A new approach for distance learning from sparse pairwise constraints," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Providence, RI, USA, Jun. 2012, pp. 2666–2672.

[5] F. Xiong, M. Gou, O. Camps, and M. Sznaier, "Person re-identification using kernel-based metric learning methods," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Zurich, Swit., Sep. 2014, pp. 1–16.

[6] W.-S. Zheng, S. Gong, and T. Xiang, "Reidentification by relative distance comparison," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 3, pp. 653–668, Mar. 2013.

[7] L. Ma, X. Yang, and D. Tao, "Person re-identification over camera networks using multi-task distance metric learning," *IEEE Trans. Image Process.*, vol. 23, no. 8, pp. 3656–3670, Aug. 2014.

[8] W. Li and X. Wang, "Locally aligned feature transforms across views," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Portland, OR, USA, Jun. 2013, pp. 3594–3601.

[9] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an ensemble of localized features," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Marseille, France, Oct. 2008, pp. 262–275.

[10] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani, "Person re-identification by symmetry-driven accumulation of local features," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, San Francisco, CA, USA, Jun. 2010, pp. 2360–2367.

[11] C. Liu, S. Gong, C. C. Loy, and X. Lin, "Person re-identification: What features are important?" in *ECCV Workshops and Demonstrations*. Springer, 2012, pp. 391–401.

[12] C.-C. Guo, S.-Z. Chen, J.-H. Lai, X.-J. Hu, and S.-C. Shi, "Multi-shot person re-identification with automatic ambiguity inference and removal," in *Proc. Int. Conf. Pattern Recognit. (ICPR)*, Stockholm, Sweden, Aug. 2014, pp. 3540–3545.

[13] R. Zhao, W. Ouyang, and X. Wang, "Unsupervised salience learning for person re-identification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Portland, OR, USA, Jun. 2013, pp. 3586–3593.

[14] ——, "Person re-identification by salience matching," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Sydney, NSW, Australia, Dec. 2013, pp. 2528–2535.

[15] M. Dikmen, E. Akbas, T. S. Huang, and N. Ahuja, "Pedestrian recognition with a learned metric," in *Proc. Asian Conf. Comput. Vis. (ACCV)*, Queenstown, New Zealand, Nov. 2010, pp. 501–512.

[16] M. Kostinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof, "Large scale metric learning from equivalence constraints," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Providence, RI, USA, Jun. 2012, pp. 2288–2295.

[17] W. Li, R. Zhao, and X. Wang, "Human reidentification with transferred metric learning," in *Proc. Asian Conf. Comput. Vis. (ACCV)*, Daejeon, Korea, Nov. 2012, pp. 31–44.

[18] Z. Li, S. Chang, F. Liang, T. S. Huang, L. Cao, and J. R. Smith, "Learning locally-adaptive decision functions for person verification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Portland, OR, USA, Jun. 2013, pp. 3610–3617.

[19] S. Pedagadi, J. Orwell, S. Velastin, and B. Boghossian, "Local Fisher discriminant analysis for pedestrian re-identification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Portland, OR, USA, Jun. 2013, pp. 3318–3325.

[20] H. Yun, P. Raman, and S. Vishwanathan, "Ranking via robust binary classification," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Montréal, Quebec, Canada, Dec. 2014, pp. 2582–2590.

[21] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Handwritten digit recognition with a back-propagation network," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*. Citeseer, 1990.

[22] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Lake Tahoe, NV, USA, Dec. 2012, pp. 1097–1105.

[23] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, Nottingham, UK, Sep. 2014.

[24] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Zurich, Swit., Sep. 2014, pp. 818–833.

[25] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 580–587.

[26] Y. Sun, Y. Chen, X. Wang, and X. Tang, "Deep learning face representation by joint identification-verification," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Montréal, Quebec, Canada, Dec. 2014, pp. 1988–1996.

[27] T. Chen, L. Lin, L. Liu, X. Luo, and X. Li, "DISC: Deep image saliency computing via progressive representation learning," *IEEE Trans. Neural Netw. Learn. Syst.*, DOI: 10.1109/TNNLS.2015.2506664, 2016.

[28] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 1653–1660.

[29] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Deep metric learning for person re-identification," in *Proc. Int. Conf. Pattern Recognit. (ICPR)*, Stockholm, Sweden, Aug. 2014, pp. 34–39.

[30] W. Li, R. Zhao, T. Xiao, and X. Wang, "DeepReID: Deep filter pairing neural network for person re-identification," in *Proc. IEEE Conf.*

*Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 152–159.

[31] S. F. Tahir and A. Cavallaro, "Cost-effective features for reidentification in camera networks," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 8, pp. 1362–1374, Aug. 2014.

[32] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino, "Custom pictorial structures for re-identification," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, vol. 2, no. 5, Dundee, UK, Aug. 2011, pp. 1–11.

[33] B. Ma, Y. Su, and F. Jurie, "BiCov: A novel image representation for person re-identification and face verification," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, Surrey, UK, Sep. 2012, pp. 57.1–57.11.

[34] I. Kviatkovsky, A. Adam, and E. Rivlin, "Color invariants for person reidentification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 7, pp. 1622–1634, Jul. 2013.

[35] Y. Yang, J. Yang, J. Yan, S. Liao, D. Yi, and S. Z. Li, "Salient color names for person re-identification," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Zurich, Swit., Sep. 2014, pp. 536–551.

[36] B. Prosser, W.-S. Zheng, S. Gong, and T. Xiang, "Person re-identification by support vector ranking," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, vol. 2, no. 5, Aberystwyth, UK, Aug. 2010, pp. 21.1–21.11.

[37] M. Sugiyama, "Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis," *J. Mach. Learn. Res.*, vol. 8, pp. 1027–1061, 2007.

[38] S. Yan, D. Xu, B. Zhang, H.-J. Zhang, Q. Yang, and S. Lin, "Graph embedding and extensions: A general framework for dimensionality reduction," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 1, pp. 40–51, Jan. 2007.

[39] R. Zhao, W. Ouyang, and X. Wang, "Learning mid-level filters for person re-identification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 144–151.

[40] J. Hu, J. Lu, and Y.-P. Tan, "Discriminative deep metric learning for face verification in the wild," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 1875–1882.

[41] P. Wu, S. C. Hoi, H. Xia, P. Zhao, D. Wang, and C. Miao, "Online multimodal deep similarity learning with application to image retrieval," in *Proc. ACM Int. Conf. Multimedia*, Barcelona, Catalunya, Spain, Oct. 2013, pp. 153–162.

[42] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 1386–1393.

[43] H. He and E. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 9, pp. 1263–1284, Sep. 2009.

[44] Y.-L. Boureau, J. Ponce, and Y. LeCun, "A theoretical analysis of feature pooling in visual recognition," in *Proc. ACM Int. Conf. on Mach. Learn. (ICML)*, Haifa, Israel, Jun. 2010, pp. 111–118.

[45] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. ACM Int. Conf. Multimedia*, Orlando, FL, USA, Nov. 2014, pp. 675–678.

[46] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *J. Mach. Learn. Res.*, vol. 10, pp. 207–244, Dec. 2009.

[47] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *Proc. ACM Int. Conf. on Mach. Learn. (ICML)*, Corvallis, OR, USA, Jun. 2007, pp. 209–216.

[48] T. Avraham, I. Gurvich, M. Lindenbaum, and S. Markovitch, "Learning implicit transfer for person re-identification," in *ECCV Workshops and Demonstrations*. Springer, 2012, pp. 381–390.

[49] N. Martinel, A. Das, C. Micheloni, and A. Roy-Chowdhury, "Re-identification in the function space of feature warps," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 8, pp. 1656–1669, Aug. 2015.

[50] Y. Wei, W. Xia, J. Huang, B. Ni, J. Dong, Y. Zhao, and S. Yan, "CNN: Single-label to multi-label," *arXiv preprint arXiv:1406.5726*, 2014.

[51] W.-S. Zheng, S. Gong, and T. Xiang, "Towards open-world person re-identification by one-shot group-based verification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 3, pp. 591–606, Mar. 2015.

[52] A. J. Ma, P. C. Yuen, and J. Li, "Domain transfer support vector ranking for person re-identification without target camera label information," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Sydney, NSW, Australia, Dec. 2013, pp. 3567–3574.

[53] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Boston, MA, USA, Jun. 2015, pp. 815–823.