# Pedestrian Alignment Network for Large-scale Person Re-identification

**Zhedong Zheng · Liang Zheng · Yi Yang**

arXiv:1707.00408v1 [cs.CV] 3 Jul 2017

**Abstract** Person re-identification (person re-ID) is mostly viewed as an image retrieval problem. This task aims to search a query person in a large image pool. In practice, person re-ID usually adopts automatic detectors to obtain cropped pedestrian images. However, this process suffers from two types of detector errors: excessive background and part missing. Both errors deteriorate the quality of pedestrian alignment and may compromise pedestrian matching due to the position and scale variances. To address the misalignment problem, we propose that alignment can be learned from an identification procedure. We introduce the pedestrian alignment network (PAN) which allows discriminative embedding learning and pedestrian alignment without extra annotations. Our key observation is that when the convolutional neural network (CNN) learns to discriminate between different identities, the learned feature maps usually exhibit strong activations on the human body rather than the background. The proposed network thus takes advantage of this attention mechanism to adaptively locate and align pedestrians within a bounding box. Visual examples show that pedestrians are better aligned with PAN. Experiments on three large-scale re-ID datasets confirm that PAN improves the discriminative ability of the feature embeddings and yields competitive accuracy with the state-of-the-art methods. [1]

Zhedong Zheng · Liang Zheng · Yi Yang
University of Technology Sydney, Australia
E-mail: zdzheng12@gmail.com

[1] The project website of this paper is `https://github.com/layumi/Pedestrian_Alignment`.



**Fig. 1** Sample images influenced by detector errors (the first row) which are aligned by the proposed method (the second row). Two types of errors are shown: excessive background and part missing. We show that the pedestrian alignment network (PAN) corrects the misalignment problem by 1) removing extra background or 2) padding zeros to the image borders. PAN reduces the scale and position variance, and the aligned output thus benefit the subsequent matching step.

## 1 Introduction

Person re-identification (person re-ID) aims at spotting the target person in different cameras, and is mostly viewed as an image retrieval problem, *i.e.,* searching for the query person in a large image pool (gallery). Recent progress mainly consists in the discriminatively learned embeddings using the convolutional neural network (CNN) on large-scale datasets. The learned embeddings extracted from the fine-tuned CNNs are shown to outperform the hand-crafted features [Zheng et al., 2016b, Xiao et al., 2016, Zhong et al., 2017].

Among the many influencing factors, misalignment is a critical one in person re-ID. This problem arises due to the usage of pedestrian detectors. In realistic set-

tings, the hand-drawn bounding boxes, existing in some previous datasets such as VIPER [Gray et al., 2007], CUHK01 [Li et al., 2012] and CUHK02 [Li and Wang, 2013], are infeasible to acquire when millions of bounding boxes are to be generated. So recent large-scale benchmarks such as CUHK03 [Li et al., 2014], Market1501 [Zheng et al., 2015] and MARS [Zheng et al., 2016a] adopt the Deformable Part Model (DPM) [Felzenszwalb et al., 2008] to automatically detect pedestrians. This pipeline largely saves the amount of labeling effort and is closer to realistic settings. However, when detectors are used, detection errors are inevitable, which may lead to two common noisy factors: excessive background and part missing. For the former, the background may take up a large proportion of a detected image. For the latter, a detected image may contain only part of the human body, *i.e.,* with missing parts (see Fig. 1).

Pedestrian alignment and re-identification are two connected problems. When we have the identity labels of the pedestrian bounding boxes, we might be able to find the optimal affine transformation that contains the most informative visual cues to discriminate between different identities. With the affine transformation, pedestrians can be better aligned. Furthermore, with superior alignment, more discriminative features can be learned, and the pedestrian matching accuracy can, in turn, be improved.

Motivated by the above-mentioned aspects, we propose to incorporate pedestrian alignment into an identification re-ID architecture, yielding the pedestrian alignment network (PAN). Given a pedestrian detected image, this network simultaneously learns to re-localize the person and categorize the person into pre-defined identities. Therefore, PAN takes advantage of the complementary nature of person alignment and re-identification.

In a nutshell, the training process of PAN is composed of the following components: 1) a network to predict the identity of an input image, 2) an affine transformation to be estimated which re-localizes the input image, and 3) another network to predict the identity of the re-localized image. For components 1) and 3), we use two convolutional branches called the base branch and alignment branch, to respectively predict the identity of the original image and the aligned image. Internally, they share the low-level features and during testing are concatenated at the fully-connected (FC) layer to generate the pedestrian descriptor. In component 2), the affine parameters are estimated using the feature maps from the high-level convolutional layer of the base branch. The affine transformation is later applied on the lower-level feature maps of the base branch. In this step, we deploy a differentiable localization network: spatial transformer network (STN) [Jaderberg et al., 2015]. With STN, we can 1) crop the detected images which may contain too much background or 2) pad zeros to the borders of images with missing parts. As a result, we reduce the impact of scale and position variances caused by misdetection and thus make pedestrian matching more precise.

Note that our method addresses the misalignment problem caused by detection errors, while the commonly used patch matching strategy aims to discover matched local structures in well-aligned images. For methods that use patch matching, it is assumed that the matched local structures locate in the same horizontal stripe [Li et al., 2014, Yi et al., 2014, Zhao et al., 2013a, Zhao et al., 2014, Liao et al., 2015, Cheng et al., 2016] or square neighborhood [Ahmed et al., 2015]. Therefore, these algorithms are robust to some small spatial variance, *e.g.,* position and scale. However, when misdetection happens, due to the limitation of the search scope, this type of methods may fail to discover the matched structures, and the risk of part mismatching may be high. Therefore, regarding the problem to be solved, the proposed method is significantly different from this line of works [Li et al., 2014, Yi et al., 2014, Ahmed et al., 2015, Zhao et al., 2013a, Zhao et al., 2014, Liao et al., 2015, Cheng et al., 2016]. We speculate that our method is a good complementary step for those using part matching.

Our contributions are summarized as follows:

– We propose the pedestrian alignment network (PAN), which simultaneously aligns pedestrians within images and learns pedestrian descriptors. Except for the identity label, we do not need any extra annotation;
– We observe that the manually cropped images are not as perfect as preassumed to be. We show that our network also improves the re-ID performance on the hand-drawn datasets which are considered to have decent person alignment.
– We achieve competitive accuracy compared to the state-of-the-art methods on three large-scale person re-ID datasets (Market-1501 [Zheng et al., 2015], CUHK03 [Li et al., 2014] and DukeMTMC-reID [Zheng et al., 2017b]).

The rest of this paper is organized as follows. Section 2 reviews and discusses related works. Section 3 illustrates the proposed method in detail. Experimental results and comparisons on three large-scale person re-ID datasets are discussed in Section 4, followed by conclusions in Section 5.

## 2 Related work

Our work aims to address two tasks: person re-identification (person re-ID) and person alignment jointly. In this section, we review the relevant works in these two domains.

### 2.1 Hand-crafted Systems for Re-ID

Person re-ID needs to find the robust and discriminative features among different cameras. Several pioneering approaches have explored person re-ID by extracting local hand-crafted features such as LBP [Mignon and Jurie, 2012], Gabor [Prosser et al., 2010] and LOMO [Liao et al., 2015]. In a series of works by [Zhao et al., 2014, Zhao et al., 2013a, Zhao et al., 2013b], the 32-dim LAB color histogram and the 128-dim SIFT descriptor are extracted from each $10 \times 10$ patches. [Zheng et al., 2015] use color name descriptor for each local patch and aggregate them into a global vector through the Bag-of-Words model. Approximate nearest neighbor search [Wang and Li, 2012] is employed for fast retrieval but accuracy compromise. [Chen et al., 2017] also deploy several different hand-crafted features extracting from overlapped body patches. Differently, [Cheng et al., 2011] localize the parts first and calculate color histograms for part-to-part correspondences. This line of works is beneficial from the local invariance in different viewpoints.

Besides finding robust feature, metric learning is nontrivial for person re-ID. [Köstinger et al., 2012] propose "KISSME" based on Mahalanobis distance and formulate the pair comparison as a log-likelihood ratio test. Further, [Liao et al., 2015] extend the Bayesian face and KISSME to learn a discriminant subspace with a metric. Aside from the methods using Mahalanobis distance, Prosser *et al.* apply a set of weak RankSVMs to assemble a strong ranker [Prosser et al., 2010]. Gray and Tao propose using the AdaBoost algorithm to fuse different features into a single similarity function [Gray and Tao, 2008]. [Loy et al., 2010] propose a cross canonical correlation analysis for the video-based person re-ID.

### 2.2 Deeply-learned Models for Re-ID

CNN-based deep learning models have been popular since [Krizhevsky et al., 2012] won ILSVRC'12 by a large margin. It extracts features and learns a classifier in an end-to-end system. More recent approaches based on CNN apply spatial constraints by splitting images or adding new patch-matching layers. [Yi et al., 2014] split a pedestrian image into three horizontal parts and respectively trained three part-CNNs to extract features. Similarly, [Cheng et al., 2016] split the convolutional map into four parts and fuse the part features with the global feature. [Li et al., 2014] add a new layer that multiplies the activation of two images in different horizontal stripes. They use this layer to allow patch matching in CNN explicitly. Later, [Ahmed et al., 2015] improve the performance by proposing a new part-matching layer that compares the activation of two images in neighboring pixels. Besides, [Varior et al., 2016a] combine CNN with some gate functions, similar to long-short-term memory (LSTM [Hochreiter and Schmidhuber, 1997]) in spirit, which aims to focus on the similar parts of input image pairs adaptively. But it is limited by the computational inefficiency because the input should be in pairs. Similarly, [Liu et al., 2016a] propose a soft attention-based model to focus on parts and combine CNN with LSTM components selectively; its limitation also consists of the computation inefficiency.

Moreover, a convolutional network has the high discriminative ability by itself without explicit patch-matching. For person re-ID, [Zheng et al., 2016c] directly use a conventional fine-tuning approach on Market-1501 [Zheng et al., 2015] and their performance outperform other recent results. [Wu et al., 2016c] combine the CNN embedding with hand-crafted features. [Xiao et al., 2016] jointly train a classification model with multiple datasets and propose a new dropout function to deal with the hundreds of identity classes. [Wu et al., 2016b] deepen the network and use filters of smaller size. [Lin et al., 2017] use person attributes as auxiliary tasks to learn more information. [Zheng et al., 2016d] propose combining the identification model with the verification model and improve the fine-tuned CNN performance. [Ding et al., 2015] and [Hermans et al., 2017] use triplet samples for training the network which considers the images from the same people and the different people at the same time. Recent work by Zheng *et al.* combined original training dataset with GAN-generated images and regularized the model [Zheng et al., 2017b]. In this paper, we adopt the similar convolutional branches without explicit part-matching layers. It is noted that we focus on a different goal on finding robust pedestrian embedding for person re-identification, and thus our method can be potentially combined with the previous methods to further improve the performance.

### 2.3 Objective Alignment

Face alignment (here refer to the rectification of face misdetection) has been widely studied. [Huang et al., 2007] propose an unsupervised method called funneled
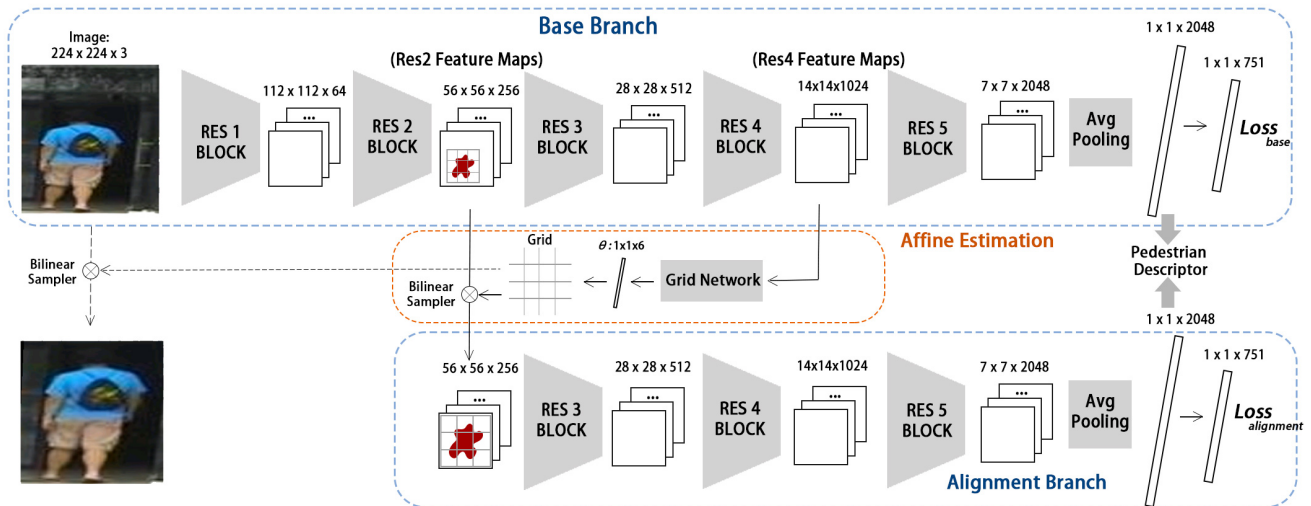
**Fig. 2** Architecture of the pedestrian alignment network (PAN). It consists of two identification networks (blue) and an affine estimation network (orange). The base branch predicts the identities from the original image. We use the high-level feature maps of the base branch (Res4 Feature Maps) to predict the grid. Then the grid is applied to the low-level feature maps (Res2 Feature Maps) to re-localize the pedestrian (red star). The alignment stream then receives the aligned feature maps to identify the person again. Note that we do not perform alignment on the original images (dotted arrow) as previously done in [Jaderberg et al., 2015] but directly on the feature maps. In the training phase, the model minimizes two identification losses. In the test phase, we concatenate two $1 \times 1 \times 2048$ FC embeddings to form a 4096-dim pedestrian descriptor for retrieval.

image to align faces according to the distribution of other images and improve this method with convolutional RBM descriptor later [Huang et al., 2012]. However, it is not trained in an end-to-end manner, and thus following tasks *i.e.,* face recognition take limited benefits from the alignment. On the other hand, several works introduce attention models for task-driven object localization. Jadeburg *et al.* [Jaderberg et al., 2015] deploy the spatial transformer network (STN) to fine-grained bird recognition and house number recognition. [Johnson et al., 2015] combine faster-RCNN [Girshick, 2015], RNN and STN to address the localization and description in image caption. Aside from using STN, Liu *et al.* use reinforcement learning to detect parts and assemble a strong model for fine-grained recognition [Liu et al., 2016c].

In person re-ID, [Baltieri et al., 2015] exploits 3D body models to the well-detected images to align the pose but does not handle the misdetection problem. Besides, the work that inspires us the most is "PoseBox" proposed by [Zheng et al., 2017a]. The PoseBox is a strengthened version of the Pictorial Structures proposed in [Cheng et al., 2011]. PoseBox is similar to our work in that 1) both works aim to solve the misalignment problem, and 2) the networks have two convolutional streams. Nevertheless, our work differs significantly from PoseBox in two aspects. First, PoseBox employs the convolutional pose machines (CPM) to generate body parts for alignment in advance, while this work learns pedestrian alignment in an end-to-end manner without extra steps. Second, PoseBox can tackle the problem of excessive background but may be less effective when some parts are missing, because CPM fails to detect body joints when the body part is absent. However, our method automatically provides solutions to both problems, *i.e.,* excessive background and part missing.

## 3 Pedestrian Alignment Network

### 3.1 Overview of PAN

Our goal is to design an architecture that jointly aligns the images and identifies the person. The primary challenge is to develop a model that supports end-to-end training and benefits from the two inter-connected tasks. The proposed architecture draws on two convolutional branches and one affine estimation branch to simultaneously address these design constraints. Fig. 2 briefly illustrates our model.

To illustrate our method, we use the ResNet-50 model [He et al., 2016] as the base model which is applied on the Market-1501 dataset [Zheng et al., 2015]. Each $Res\_i, i = 1, 2, 3, 4, 5$ block in Fig. 2 denotes several convolutional layers with batch normalization, ReLU, and optionally max pooling. After each block, the feature
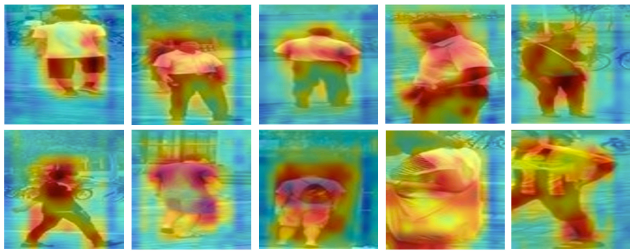
**Fig. 3** We visualize the Res4 Feature Maps in the base branch. We observe that high responses are mostly concentrated on the pedestrian body. So we use the Res4 Feature Maps to estimate the affine parameters.

maps are down-sampled to be half of the size of the feature maps in the previous block. For example, $Res\_1$ down-samples the width and height of an image from $224 \times 224$ to $112 \times 112$. In Section 3.2 and Section 3.3, we first describe the convolutional branches and affine estimation branches of our model. Then in Section 3.4 we address the details of a pedestrian descriptor. When testing, we use the descriptor to retrieve the query person. Further, we discuss the re-ranking method as a subsequent processing in Section 3.5.

### 3.2 Base and Alignment Branches

Recent progress in person re-ID datasets allows the CNN model to learn more discriminative visual representations. There are two main convolutional branches exist in our model, called the base branch and the alignment branch. Both branches are classification networks that predict the identity of the training images. Given an originally detected image, the base branch not only learns to distinguish its identity from the others but also encodes the appearance of the detected image and provides the clues for the spatial localization (see Fig. 3). The alignment branch shares a similar convolutional network but processes the aligned feature maps produced by the affine estimation branch.

In the base branch, we train the ResNet-50 model [He et al., 2016], which consists of five down-sampling blocks and one global average pooling. We deploy the model pre-trained on ImageNet [Deng et al., 2009] and remove the final fully-connected (FC) layer. There are $K = 751$ identities in the Market-1501 training set, so we add an FC layer to map the CNN embedding of size $1 \times 1 \times 2048$ to 751 unnormalized probabilities. The alignment branch, on the other hand, is comprised of three ResBlocks and one average pooling layer. We also add an FC layer to predict the multi-class probabilities. The two branches do not share weight. We use $W_1$ and

$W_2$ to denote the parameters of the two convolutional branches, respectively.

More formally, given an input image $x$, we use $p(k|x)$ to denote the probability that the image $x$ belongs to the class $k \in \{1...K\}$. Specifically, $p(k|x) = \frac{exp(z_k)}{\sum_{k=1}^{K} exp(z_i)}$. Here $z_i$ is the outputted probability from the CNN model. For the two branches, the cross-entropy losses are formulated as:

$$l_{base}(W_1, x, y) = -\sum_{k=1}^{K} (log(p(k|x))q(k|x)), \qquad (1)$$

$$l_{align}(W_2, x_a, y) = -\sum_{k=1}^{K} (log(p(k|x_a))q(k|x_a)), \qquad (2)$$

where $x_a$ denotes the aligned input. It can be derived from the original input $x_a = T(x)$. Given the label $y$, the ground-truth distribution $q(y|x) = 1$ and $q(k|x) = 0$ for all $k \neq y$. If we discard the 0 term in Eq. 1 and Eq. 2, the losses are equivalent to:

$$l_{base}(W_1, x, y) = -log(p(y|x)), \qquad (3)$$

$$l_{align}(W_2, x_a, y) = -log(p(y|x_a)). \qquad (4)$$

Thus, at each iteration, we wish to minimize the total entropy, which equals to maximizing the possibility of the correct prediction.

### 3.3 Affine Estimation Branch

To address the problems of excessive background and part missing, the key idea is to predict the position of the pedestrian and do the corresponding spatial transform. When excessive background exists, a cropping strategy should be used; under part missing, we need to pad zeros to the corresponding image borders. Both strategies need to find the parameters for the affine transformation. In this paper, this function is implemented by the affine estimation branch.

The affine estimation branch receives two input tensors of activations $14 \times 14 \times 1024$ and $56 \times 56 \times 256$ from the base branch. We name the two tensors the Res2 Feature Maps and the Res4 Feature Maps, respectively. The Res4 Feature Maps contain shallow feature maps of the original image and reflects the local pattern information. On the other hand, since the Res2 Feature Maps are closer to the classification layer, it encodes the attention on the pedestrian and semantic cues for aiding identification. The affine estimation branch contains one bilinear sampler and one small network called

Grid Network. The Grid Network contains one Res-Block and one average pooling layer. We pass the Res4 Feature Maps through Grid Network to regress a set of 6-dimensional transformer parameters. The learned transforming parameters $\theta$ are used to produce the image grid. The affine transformation process can be written as below,

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} \ \theta_{12} \ \theta_{13} \\ \theta_{21} \ \theta_{22} \ \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}, \tag{5}$$

where $(x_i^t, y_i^t)$ are the target coordinates on the output feature map, and $(x_i^s, y_i^s)$ are the source coordinates on the input feature maps (Res2 Feature Maps). $\theta_{11}, \theta_{12}, \theta_{21}$ and $\theta_{22}$ deal with the scale and rotation transformation, while $\theta_{13}$ and $\theta_{23}$ deal with the offset. In this paper, we set the coordinates as: (-1,-1) refer to the pixel on the top left of the image, while (1,1) refer to the bottom right pixel. For example, if $\theta = \begin{bmatrix} 0.8 \ \ 0 \ \ -0.1 \\ 0 \ \ 0.7 \ \ 0 \end{bmatrix}$, the pixel value of point $(-1, -1)$ on the output image is equal to that of $(-0.9, -0.7)$ on the original map. We use a bilinear sampler to make up the missing pixels, and we assign zeros to the pixels located out of the original range. So we obtain an injective function from the original feature map $V$ to the aligned output $U$. More formally, we can formulate the equation:

$$U_{(m,n)}^c =$$
$$\sum_{x^s}^{H} \sum_{y^s}^{W} V_{(x^s, y^s)}^c max(0, 1-|x^t-m|)max(0, 1-|y^t-n|). \tag{6}$$

$U_{(m,n)}^c$ is the output feature map at location $(m, n)$ in channel $c$, and $V_{(x^s, y^s)}^c$ is the input feature map at location $(x_s, y_s)$ in channel $c$. If $(x_t, y_t)$ is close to $(m, n)$, we add the pixel at $(x_s, y_s)$ according to the bilinear sampling.

In this work, we do not perform pedestrian alignment on the original image; instead, we choose to achieve an equivalent function on the shallow feature maps. By using the feature maps, we reduce the running time and parameters of the model. This explains why we apply re-localization grid on the feature maps. The bilinear sampler receives the grid, and the feature maps to produce the aligned output $x_a$. We provide some visualization examples in Fig. 3. Res4 Feature maps are shown. We can observe that through ID supervision, we are able to re-localize the pedestrian and correct misdetections to some extent.

### 3.4 Pedestrian Descriptor

Given the fine-tuned PAN model and an input image $x_i$, the pedestrian descriptor is the weighted fusion of the FC features of the base branch and the alignment branch. That is, we are able to capture the pedestrian characteristic from the original image and the aligned image. In the Section 4.3, the experiment shows that the two features are complementary to each other and improve the person re-ID performance.

In this paper, we adopt a straightforward late fusion strategy, i.e., $f_i = g(f_i^1, f_i^2)$. Here $f_i^1$ and $f_i^2$ are the FC descriptors from two types of images, respectively. We reshape the tensor after the final average pooling to a 1-dim vector as the pedestrian descriptor of each branch. The pedestrian descriptor is then written as:

$$f_i = \left[ \alpha |f_i^1|^{\mathrm{T}}, (1-\alpha)|f_i^2|^{\mathrm{T}} \right]^{\mathrm{T}}. \tag{7}$$

The $|\cdot|$ operator denotes an $l^2$-normalization step. We concatenate the aligned descriptor with the original descriptor, both after $l^2$-normalization. $\alpha$ is the weight for the two descriptors. If not specified, we simply use $\alpha = 0.5$ in our experiments.

### 3.5 Re-ranking for re-ID

In this work, we first obtain the rank list $N(q, n) = [x_1, x_2, ...x_n]$ by sorting the Euclidean distance of gallery images to the query. Distance is calculated as $D_{i,j} = (f_i - f_j)^2$, where $f_i$ and $f_j$ are $l_2$-normalized features of image $i$ and $j$, respectively. We then perform re-ranking to obtain better retrieval results. Several re-ranking methods can be applied in person re-ID [Ye et al., 2015, Qin et al., 2011, Zhong et al., 2017]. Specifically, we follow a state-of-the-art re-ranking method proposed in [Zhong et al., 2017].

Apart from the Euclidean distance, we additionally consider the Jaccard similarity. To introduce this distance, we first define a robust retrieval set for each image. The k-reciprocal nearest neighbors $R(p, k)$ are comprised of such element that appears in the top-k retrieval rank of the query $p$ while the query is in the top-k rank of the element as well.

$$R(p, k) = \{x | x \in N(p, k), p \in N(x, k)\} \tag{8}$$

According to [Zhong et al., 2017], we extend the set $R$ to $R^*$ to include more positive samples. Taking the advantage of set $R^*$, we use the Jaccard similarity for re-ranking. When we use the correctly matched images to conduct the retrieval, we should retrieve a similar

**Fig. 4** Sample images from Market-1501 [Zheng et al., 2015], CUHK03 [Li et al., 2014] and DukeMTMC-reID [Zheng et al., 2017b]. The three datasets are collected in different scenes with different detection bias.

rank list as we use the original query. The Jaccard distance can be simply formulated as:

$$D_{similarity} = 1 - \frac{|R^*(q,k) \cap R^*(x_i,k)|}{|R^*(q,k) \cup R^*(x_i,k)|}, \tag{9}$$

where $|\cdot|$ denotes the cardinality of the set. If $R^*(q,k)$ and $R^*(x_i,k)$ share more elements, $x_i$ is more likely to be a true match. This usually helps us to distinguish some hard negative samples from the correct matches. During testing, this similarity distance is added to the Euclidean distance to re-rank the result. In the experiment, we show that re-ranking further improves our results.

## 4 Experiments

In this section, we report the results on three large-scale datasets: Market-1501 [Zheng et al., 2015], CUHK03 [Li et al., 2014] and DukeMTMC-reID [Zheng et al., 2017b]. As for the detector, Market-1501 and CUHK03 (detected) datasets are automatically detected by DPM and face the misdetection problem. It is unknown if the manually annotated images after slight alignment would bring any extra benefit. So we also evaluate the proposed method on the manually annotated images of CUHK03 (labeled) and DukeMTMC-reID, which consist of hand-drawn bounding boxes. As shown in Fig. 4, the three datasets have different characteristics, *i.e.*, scene variances, and detection bias.

### 4.1 Datasets

**Market1501** is a large-scale person re-ID dataset collected in a university campus. It contains 19,732 gallery images, 3,368 query images and 12,936 training images collected from six cameras. There are 751 identities in the training set and 750 identities in the testing set without overlapping. Every identity in the training set has 17.2 photos on average. All images are automatically detected by the DPM detector [Felzenszwalb et al., 2008]. The misalignment problem is common, and the dataset is close to the realistic settings. We use all the 12,936 detected images to train the network and follow the evaluation protocol in the original dataset. There are two evaluation settings. A single query is to use one image of one person as query, and multiple query means to use several images of one person under one camera as a query.

**CUHK03** contains 14,097 images of 1,467 identities [Li et al., 2014]. Each identity contains 9.6 images on average. We follow the new training/testing protocol proposed in [Zhong et al., 2017] to evaluate the multi-shot re-ID performance. This setting uses a larger testing gallery and is different from the papers published earlier than [Zhong et al., 2017], such as [Liu et al., 2016a] and [Varior et al., 2016b]. There are 767 identities in the training set and 700 identities in the testing set (The former setting uses 1,367 IDs for training and the other 100 IDs for testing). Since we usually face a large-scale searching image pool cropped from surveillance videos, a larger testing pool is closer to the realistic image retrieval setting. In the "detected" set, all the bounding boxes are produced by DPM; in the "labeled" set, the images are all hand-drawn. In this paper, we evaluate our method on "detected" and "labeled" sets, respectively. If not specified, "CUHK03" denotes the detected set, which is more challenging.

**DukeMTMC-reID** is a subset of the DukeMTMC [Ristani et al., 2016] and contains 36,411 images of 1,812 identities shot by eight high-resolution cameras. It is one of the largest pedestrian image datasets. The pedestrian images are cropped from hand-drawn bounding boxes. The dataset consists 16,522 training images of 702 identities, 2,228 query images of the other 702 identities and 17,661 gallery images. It is challenging because many pedestrians are in the similar clothes, and may be occluded by cars or trees. We follow the evaluation protocol in [Zheng et al., 2017b].

**Evaluation Metrics.** We evaluate our method with rank-1, 5, 20 accuracy and mean average precision (mAP). Here, rank-$i$ accuracy denotes the probability whether one or more correctly matched images appear in top-$i$. If no correctly matched images appear in the top-$i$

| Methods | dim | Market-1501 | | | | CUHK03 (detected) | | | | CUHK03 (labeled) | | | | DukeMTMC-reID | | | |
|---------|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| | | 1 | 5 | 20 | mAP | 1 | 5 | 20 | mAP | 1 | 5 | 20 | mAP | 1 | 5 | 20 | mAP |
| Base | 2,048 | 80.17 | 91.69 | 96.59 | 59.14 | 30.50 | 51.07 | 71.64 | 29.04 | 31.14 | 52.00 | 74.21 | 29.80 | 65.22 | 79.13 | 87.75 | 44.99 |
| Alignment | 2,048 | 79.01 | 90.86 | 96.14 | 58.27 | 34.14 | 54.50 | 72.71 | 31.71 | 35.29 | 53.64 | 72.43 | 32.90 | 68.36 | 81.37 | 88.64 | 47.14 |
| PAN | 4,096 | 82.81 | 93.53 | 97.06 | 63.35 | 36.29 | 55.50 | 75.07 | 34.00 | 36.86 | 56.86 | 75.14 | 35.03 | 71.59 | 83.89 | 90.62 | 51.51 |

**Table 1** Comparison of different methods on Market-1501, CUHK03 (detected), CUHK03 (labeled) and DukeMTMC-reID. Rank-1, 5, 20 accuracy (%) and mAP (%) are shown. Note that the base branch is the same as the classification baseline [Zheng et al., 2016b]. We observe consistent improvement of our method over the individual branches on the three datasets.

of the retrieval list, rank-$i = 0$, otherwise rank-$i = 1$. We report the mean rank-$i$ accuracy for query images. On the other hand, for each query, we calculate the area under the Precision-Recall curve, which is known as the average precision (AP). The mean of the average precision (mAP) then is calculated, which reflects the precision and recall rate of the performance.

## 4.2 Implementation Details

**ConvNet.** In this work, we first fine-tune the base branch on the person re-ID datasets. Then, the base branch is fixed while we fine-tune the whole network. Specifically, when fine-tuning the base branch, the learning rate decrease from $10^{-3}$ to $10^{-4}$ after 30 epochs. We stop training at the 40th epoch. Similarly, when we train the whole model, the learning rate decrease from $10^{-3}$ to $10^{-4}$ after 30 epochs. We stop training at the 40th epoch. We use mini-batch stochastic gradient descent with a Nesterov momentum fixed to 0.9 to update the weights. Our implementation is based on the Matconvnet [Vedaldi and Lenc, 2015] package. The input images are uniformly resized to the size of $224 \times 224$. In addition, we perform simple data augmentation such as cropping and horizontal flipping following [Zheng et al., 2016d].

**STN.** For the affine estimation branch, the network may fall into a local minimum in early iterations. To stabilize training, we find that a small learning rate is useful. We, therefore, use a learning rate of $1 \times 10^{-5}$ for the final convolutional layer in the affine estimation branch. In addition, we set the all $\theta = 0$ except that $\theta_{11}, \theta_{22} = 0.8$ and thus, the alignment branch starts training from looking at the center part of the Res2 Feature Maps.

## 4.3 Evaluation

**Evaluation of the ResNet baseline.** We implement the baseline according to the routine proposed in [Zheng et al., 2016b], with the details specified in Section 4.2.



**Fig. 5** Sensitivity of person re-ID accuracy to parameter $\alpha$. Rank-1 accuracy(%) and mAP(%) on three datasets are shown.

We report our baseline results in Table 1. The rank-1 accuracy is 80.17%, 30.50%, 31.14% and 65.22% on Market1501, CUHK03 (detected), CUHK03 (labeled) and DukeMTMC-reID respectively. The baseline model is on par with the network in [Zheng et al., 2016d, Zheng et al., 2016b]. In our recent implementation, we use a smaller batch size of 16 and a dropout rate of 0.75. We have therefore obtained a higher baseline rank-1 accuracy 80.17% on Market-1501 than 73.69% in the [Zheng et al., 2016d, Zheng et al., 2016b]. For a fair comparison, we will present the results of our methods built on this new baseline. Note that this baseline result itself is higher than many previous works [Barbosa et al., 2017, Varior et al., 2016a, Zheng et al., 2017a, Zheng et al., 2016d].

**Base branch. vs. alignment branch** To investigate how alignment helps to learn discriminative pedestrian representations, we evaluate the Pedestrian descriptors of the base branch and the alignment branch, respectively. Two conclusions can be inferred.

First, as shown in Table 1, the alignment branch yields higher performance *i.e.*, +3.64% / +4.15% on the two dataset settings (CUHK03 detected/labeled) and +3.14% on DukeMTMC-reID, and achieves a very similar result with the base branch on Market-1501. We speculate that Market-1501 contains more intensive detection errors than the other three datasets and thus, the effect of alignment is limited.

Second, although the CUHK (labeled) dataset and the DukeMTMC-reID dataset are manually annotated,

| Method | Single Query | | Multi. Query | |
| --- | --- | --- | --- | --- |
| | rank-1 | mAP | rank-1 | mAP |
| DADM [Su et al., 2016] | 39.4 | 19.6 | 49.0 | 25.8 |
| BoW+kissme [Zheng et al., 2015] | 44.42 | 20.76 | - | - |
| MR-CNN [Ustinova et al., 2015]* | 45.58 | 26.11 | 56.59 | 32.26 |
| MST-CNN [Liu et al., 2016b] | 45.1 | - | 55.4 | - |
| FisherNet [Wu et al., 2016a]* | 48.15 | 29.94 | - | - |
| CAN [Liu et al., 2016a]* | 48.24 | 24.43 | - | - |
| SL [Chen et al., 2016] | 51.90 | 26.35 | - | - |
| S-LSTM [Varior et al., 2016b] | - | - | 61.6 | 35.3 |
| DNS [Zhang et al., 2016] | 55.43 | 29.87 | 71.56 | 46.03 |
| Gate Reid [Varior et al., 2016a] | 65.88 | 39.55 | 76.04 | 48.45 |
| SOMAnet [Barbosa et al., 2017]* | 73.87 | 47.89 | 81.29 | 56.98 |
| PIE [Zheng et al., 2017a]* | 78.65 | 53.87 | - | - |
| Verif.-Identif. [Zheng et al., 2016d]* | 79.51 | 59.87 | 85.84 | 70.33 |
| SVDNet [Sun et al., 2017]* | 82.3 | 62.1 | - | - |
| DeepTransfer [Geng et al., 2016]* | 83.7 | 65.5 | 89.6 | 73.8 |
| GAN [Zheng et al., 2017b]* | 83.97 | 66.07 | 88.42 | 76.10 |
| APR [Lin et al., 2017]* | 84.29 | 64.67 | - | - |
| Triplet [Hermans et al., 2017]* | 84.92 | 69.14 | 90.53 | 76.42 |
| Triplet+re-rank [Hermans et al., 2017]* | 86.67 | 81.07 | **91.75** | 87.18 |
| Basel. | 80.17 | 59.14 | 87.41 | 72.05 |
| Ours | 82.81 | 63.35 | 88.18 | 71.72 |
| Ours+re-rank | 85.78 | 76.56 | 89.79 | 83.79 |
| Ours (GAN) | 86.67 | 69.33 | 90.88 | 76.32 |
| Ours (GAN)+re-rank | **88.57** | **81.53** | 91.45 | **87.44** |

**Table 2** Rank-1 precision (%) and mAP (%) on Market-1501. We also provide results of the fine-tuned ResNet50 baseline which has the same accuracy with the base branch. * the respective paper is on ArXiv but not published.

the alignment branch still improves the performance of the base branch. This observation demonstrates that the manual annotations may not be good enough for the machine to learn a good descriptor. In this scenario, alignment is non-trivial and makes the pedestrian representation more discriminative.

**The complementary of the two branches.** As mentioned, the two descriptors capture the different pedestrian characteristic from the original image and the aligned image. We follow the setting in Section 3.4 and simply combine the two features to form a stronger pedestrian descriptor. The results are summarized in Table 1. We observe a constant improvement on the three datasets when we concatenate the two branch descriptors. The fused descriptor improves +2.64%, +2.15%, +1.63% and 3.23% on Market-1501, CUHK03(detected), CUHK03(labeled) and DukeMTMC-reID, respectively. The two branches are complementary and thus, contain more meaningful information than a separate branch. Aside from the improvement of the accuracy, this simple fusion is efficient sine it does not introduce additional computation.

**Parameter sensitivity.** We evaluate the sensitivity of the person re-ID accuracy to the parameter $\alpha$. As shown in Fig. 5, we report the rank-1 accuracy and mAP when tuning the $\alpha$ from 0 to 1. We observe that the change of rank-1 accuracy and mAP are relatively

| Method | rank-1 | mAP |
| --- | --- | --- |
| BoW+kissme [Zheng et al., 2015] | 25.13 | 12.17 |
| LOMO+XQDA [Liao et al., 2015] | 30.75 | 17.04 |
| Gan [Zheng et al., 2017b] | 67.68 | 47.13 |
| OIM [Xiao et al., 2017] | 68.1 | - |
| APR [Lin et al., 2017] | 70.69 | 51.88 |
| SVDNet [Sun et al., 2017] | **76.7** | 56.8 |
| Basel. [Zheng et al., 2017b] | 65.22 | 44.99 |
| Ours | 71.59 | 51.51 |
| Ours + re-rank | 75.94 | **66.74** |

**Table 3** Rank-1 accuracy (%) and mAP (%) on DukeMTMC-reID. We follow the evaluation protocol in [Zheng et al., 2017b]. We also provide the result of the fine-tuned ResNet50 baseline for fair comparison.
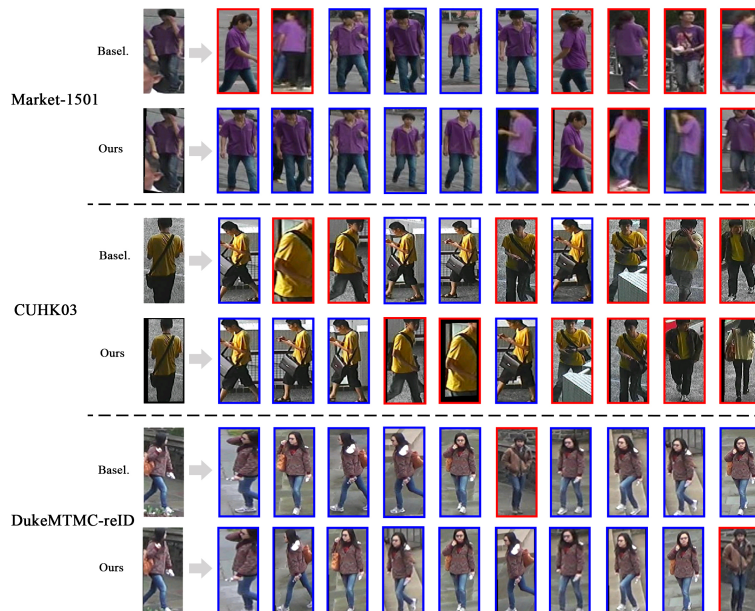
small corresponding to the $\alpha$. Our reported result simply use $\alpha = 0.5$. $\alpha = 0.5$ may not be the best choice for a particular dataset. But if we do not foreknow the distribution of the dataset, it is a simple and straightforward choice.

**Comparison with the state-of-the-art methods.** We compare our method with the state-of-the-art methods on Market-1501, CUHK03 and DukeMTMC-reID in Table 2, Table 4 and Table 3, respectively. On Market-1501, we achieve **rank-1 accuracy = 85.78%, mAP = 76.56%** after re-ranking, which is the best result compared to the published paper, and the second

| Method | Detected | | Labeled | |
| --- | --- | --- | --- | --- |
| | rank-1 | mAP | rank-1 | mAP |
| BoW+XQDA [Zheng et al., 2015] | 6.36 | 6.39 | 7.93 | 7.29 |
| LOMO+XQDA [Liao et al., 2015] | 12.8 | 11.5 | 14.8 | 13.6 |
| ResNet50+XQDA [Zhong et al., 2017] | 31.1 | 28.2 | 32.0 | 29.6 |
| ResNet50+XQDA+re-rank [Zhong et al., 2017] | 34.7 | 37.4 | 38.1 | 40.3 |
| Basel. | 30.5 | 29.0 | 31.1 | 29.8 |
| Ours | 36.3 | 34.0 | 36.9 | 35.0 |
| Ours+re-rank | **41.9** | **43.8** | **43.9** | **45.8** |

**Table 4** Rank-1 accuracy (%) and mAP (%) on CUHK03 using the new evaluation protocol in [Zhong et al., 2017]. This setting uses a larger testing gallery and is different from the papers published earlier than [Zhong et al., 2017], such as [Liu et al., 2016a] and [Varior et al., 2016b]. There are 767 identities in the training set and 700 identities in the testing set (The former setting uses 1,367 IDs for training and the other 100 IDs for testing). Since we usually face a large-scale searching image pool cropped from surveillance videos, a larger testing pool is more challenging and closer to the realistic image retrieval setting. So we evaluate the proposed method on the "detected" and "labeled" subsets according to this new multi-shot protocol. We also provide the result of our fine-tuned ResNet50 baseline for fair comparison.

**Fig. 6** Sample retrieval results on the three datasets. The images in the first column are queries. The retrieved images are sorted according to the similarity score from left to right. For each query, the first row shows the result of baseline [Zheng et al., 2016b], and the second row denotes the results of PAN. The correct and false matches are in the blue and red rectangles, respectively. Images in the rank lists obtained by PAN demonstrate amelioration in alignment. Best viewed when zoomed in.



best among all the available results including the arXiv paper. Our model is also adaptive to previous models. One of the previous best results is based on the model regularized by GAN [Zheng et al., 2017b]. We combine the model trained on GAN generated images and thus, achieve the state-of-the-art result **rank-1 accuracy = 88.57%, mAP = 81.53%** on Market-1501. On CUHK03, we arrive at a competitive result **rank-1 accuracy = 36.3%, mAP=34.0%** on the detected dataset and **rank-1 accuracy = 36.9%, mAP = 35.0%** on the labeled dataset. After re-ranking, we further achieve a state-of-the art result **rank-1 accuracy = 41.9%, mAP=43.8%** on the detected dataset and **rank-1 accuracy = 43.9%, mAP = 45.8%** on the labeled dataset. On DukeMTMC-reID, we also observe a state-of-the-art result **rank-1 accuracy = 75.94% and mAP = 66.74%** after re-ranking. Despite the

visual disparities among the three datasets, *i.e.,* scene variance, and detection bias, we show that our method consistently improves the re-ID performance.

As shown in Fig. 6, we visualize some retrieval results on the three datasets. Images in the rank lists obtained by PAN demonstrate amelioration in alignment. Comparing to the baseline, true matches which are misaligned originally receive higher ranks, while false matches have lower ranks

**Visualization of the alignment.** We further visualize the aligned images in Fig. 7. As aforementioned, the proposed network does not process the alignment on the original image. To visualize the aligned images, we extract the predicted affine parameters and then apply the affine transformation on the originally detected images manually. **We observe that the network does not perform perfect alignment as the human,**

**Fig. 7** Examples of pedestrian images before and after alignment on three datasets (Market-1501, DukeMTMC-reID and CUHK03). Pairs of input images and aligned images are shown. By removing excessive background or padding zeros to image borders, we observe that PAN reduces the scale and location variance.

**but it more or less reduces the scale and position variance, which is critical for the network to learn the representations.** So the proposed network improves the performance of the person re-ID.

## 5 Conclusion

Pedestrian alignment and re-identification are two interconnected problems, which inspires us to develop an attention-based system. In this work, we propose the pedestrian alignment network (PAN), which simultaneously aligns the pedestrians within bounding boxes and learns the pedestrian descriptors. Taking advantage of the attention of CNN feature maps to the human body, PAN addresses the misalignment problem and person re-ID together and thus, improves the person re-ID accuracy. Except for the identity label, we do not need any extra annotation. We also observe that the manually cropped images are not as perfect as preassumed to be. Our network also improves the re-ID performance on the datasets with hand-drawn bounding boxes. Experiments on three different datasets indicate that our method is competitive with the state-of-the-art methods. In the future, we will continue to investigate the

attention-based model and apply our model to other fields *i.e.,* car recognition.

## References

Ahmed et al., 2015. Ahmed, E., Jones, M., and Marks, T. K. (2015). An improved deep learning architecture for person re-identification. In *CVPR*.

Baltieri et al., 2015. Baltieri, D., Vezzani, R., and Cucchiara, R. (2015). Mapping appearance descriptors on 3d body models for people re-identification. *IJCV*.

Barbosa et al., 2017. Barbosa, I. B., Cristani, M., Caputo, B., Rognhaugen, A., and Theoharis, T. (2017). Looking beyond appearances: Synthetic training data for deep cnns in re-identification. *arXiv:1701.03153*.

Chen et al., 2016. Chen, D., Yuan, Z., Chen, B., and Zheng, N. (2016). Similarity learning with spatial constraints for person re-identification. In *CVPR*.

Chen et al., 2017. Chen, D., Yuan, Z., Wang, J., Chen, B., Hua, G., and Zheng, N. (2017). Exemplar-guided similarity learning on polynomial kernel feature map for person re-identification. *IJCV*.

Cheng et al., 2016. Cheng, D., Gong, Y., Zhou, S., Wang, J., and Zheng, N. (2016). Person re-identification by multi-channel parts-based cnn with improved triplet loss function. In *CVPR*.

Cheng et al., 2011. Cheng, D. S., Cristani, M., Stoppa, M., Bazzani, L., and Murino, V. (2011). Custom pictorial structures for re-identification. In *BMVC*.

Deng et al., 2009. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *CVPR*.

Ding et al., 2015. Ding, S., Lin, L., Wang, G., and Chao, H. (2015). Deep feature learning with relative distance comparison for person re-identification. *Pattern Recognition*, 48(10):2993–3003.

Felzenszwalb et al., 2008. Felzenszwalb, P., McAllester, D., and Ramanan, D. (2008). A discriminatively trained, multiscale, deformable part model. In *CVPR*.

Geng et al., 2016. Geng, M., Wang, Y., Xiang, T., and Tian, Y. (2016). Deep transfer learning for person re-identification. *arXiv:1603.06765*.

Girshick, 2015. Girshick, R. (2015). Fast r-cnn. In *ICCV*.

Gray et al., 2007. Gray, D., Brennan, S., and Tao, H. (2007). Evaluating appearance models for recognition, reacquisition, and tracking. In *PETS*, volume 3. Citeseer.

Gray and Tao, 2008. Gray, D. and Tao, H. (2008). Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*.

He et al., 2016. He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *CVPR*.

Hermans et al., 2017. Hermans, A., Beyer, L., and Leibe, B. (2017). In defense of the triplet loss for person re-identification. *arXiv:1703.07737*.

Hochreiter and Schmidhuber, 1997. Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

Huang et al., 2012. Huang, G., Mattar, M., Lee, H., and Learned-Miller, E. G. (2012). Learning to align from scratch. In *NIPS*.

Huang et al., 2007. Huang, G. B., Jain, V., and Learned-Miller, E. (2007). Unsupervised joint alignment of complex images. In *ICCV*.

Jaderberg et al., 2015. Jaderberg, M., Simonyan, K., Zisserman, A., et al. (2015). Spatial transformer networks. In *NIPS*.

Johnson et al., 2015. Johnson, J., Karpathy, A., and Fei-Fei, L. (2015). Densecap: Fully convolutional localization networks for dense captioning. *arXiv:1511.07571*.

Köstinger et al., 2012. Köstinger, M., Hirzer, M., Wohlhart, P., Roth, P. M., and Bischof, H. (2012). Large scale metric learning from equivalence constraints. In *CVPR*.

Krizhevsky et al., 2012. Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *NIPS*.

Li and Wang, 2013. Li, W. and Wang, X. (2013). Locally aligned feature transforms across views. In *CVPR*.

Li et al., 2012. Li, W., Zhao, R., and Wang, X. (2012). Human reidentification with transferred metric learning. In *ACCV*.

Li et al., 2014. Li, W., Zhao, R., Xiao, T., and Wang, X. (2014). Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*.

Liao et al., 2015. Liao, S., Hu, Y., Zhu, X., and Li, S. Z. (2015). Person re-identification by local maximal occurrence representation and metric learning. In *CVPR*.

Lin et al., 2017. Lin, Y., Zheng, L., Zheng, Z., Wu, Y., and Yang, Y. (2017). Improving person re-identification by attribute and identity learning. *arXiv:1703.07220*.

Liu et al., 2016a. Liu, H., Feng, J., Qi, M., Jiang, J., and Yan, S. (2016a). End-to-end comparative attention networks for person re-identification. *arXiv:1606.04404*.

Liu et al., 2016b. Liu, J., Zha, Z.-J., Tian, Q., Liu, D., Yao, T., Ling, Q., and Mei, T. (2016b). Multi-scale triplet cnn for person re-identification. In *ACM Multimedia*.

Liu et al., 2016c. Liu, X., Xia, T., Wang, J., Yang, Y., Zhou, F., and Lin, Y. (2016c). Fully convolutional attention networks for fine-grained recognition. *arXiv:1611.05244*.

Loy et al., 2010. Loy, C. C., Xiang, T., and Gong, S. (2010). Time-delayed correlation analysis for multi-camera activity understanding. *IJCV*.

Mignon and Jurie, 2012. Mignon, A. and Jurie, F. (2012). Pcca: A new approach for distance learning from sparse pairwise constraints. In *CVPR*.

Prosser et al., 2010. Prosser, B., Zheng, W.-S., Gong, S., Xiang, T., and Mary, Q. (2010). Person re-identification by support vector ranking. In *BMVC*.

Qin et al., 2011. Qin, D., Gammeter, S., Bossard, L., Quack, T., and Van Gool, L. (2011). Hello neighbor: Accurate object retrieval with k-reciprocal nearest neighbors. In *CVPR*.

Ristani et al., 2016. Ristani, E., Solera, F., Zou, R., Cucchiara, R., and Tomasi, C. (2016). Performance measures and a data set for multi-target, multi-camera tracking. In *ECCVW*.

Su et al., 2016. Su, C., Zhang, S., Xing, J., Gao, W., and Tian, Q. (2016). Deep attributes driven multi-camera person re-identification. *ECCV*.

Sun et al., 2017. Sun, Y., Zheng, L., Deng, W., and Wang, S. (2017). Svdnet for pedestrian retrieval. *arXiv:1703.05693*.

Ustinova et al., 2015. Ustinova, E., Ganin, Y., and Lempitsky, V. (2015). Multiregion bilinear convolutional neural networks for person re-identification. *arXiv:1512.05300*.

Varior et al., 2016a. Varior, R. R., Haloi, M., and Wang, G. (2016a). Gated siamese convolutional neural network architecture for human re-identification. In *ECCV*.

Varior et al., 2016b. Varior, R. R., Shuai, B., Lu, J., Xu, D., and Wang, G. (2016b). A siamese long short-term memory architecture for human re-identification. In *ECCV*.

Vedaldi and Lenc, 2015. Vedaldi, A. and Lenc, K. (2015). Matconvnet – convolutional neural networks for matlab. In *ACM Multimedia*.

Wang and Li, 2012. Wang, J. and Li, S. (2012). Query-driven iterated neighborhood graph search for large scale indexing. In *ACM Multimedia*.

Wu et al., 2016a. Wu, L., Shen, C., and Hengel, A. v. d. (2016a). Deep linear discriminant analysis on fisher networks: A hybrid architecture for person re-identification. *arXiv:1606.01595*.

Wu et al., 2016b. Wu, L., Shen, C., and Hengel, A. v. d. (2016b). Personnet: Person re-identification with deep convolutional neural networks. *arXiv:1601.07255*.

Wu et al., 2016c. Wu, S., Chen, Y.-C., Li, X., Wu, A.-C., You, J.-J., and Zheng, W.-S. (2016c). An enhanced deep feature representation for person re-identification. In *WACV*.

Xiao et al., 2016. Xiao, T., Li, H., Ouyang, W., and Wang, X. (2016). Learning deep feature representations with domain guided dropout for person re-identification. In *CVPR*.

Xiao et al., 2017. Xiao, T., Li, S., Wang, B., Lin, L., and Wang, X. (2017). Joint detection and identification feature learning for person search. *arXiv preprint arXiv:1604.01850*.

Ye et al., 2015. Ye, M., Liang, C., Wang, Z., Leng, Q., and Chen, J. (2015). Ranking optimization for person re-identification via similarity and dissimilarity. In *ACM Multimedia*.

Yi et al., 2014. Yi, D., Lei, Z., Liao, S., and Li, S. Z. (2014). Deep metric learning for person re-identification. In *ICPR*.

Zhang et al., 2016. Zhang, L., Xiang, T., and Gong, S. (2016). Learning a discriminative null space for person re-identification. *CVPR*.

Zhao et al., 2013a. Zhao, R., Ouyang, W., and Wang, X. (2013a). Person re-identification by salience matching. In *ICCV*.

Zhao et al., 2013b. Zhao, R., Ouyang, W., and Wang, X. (2013b). Unsupervised salience learning for person re-identification. In *CVPR*.

Zhao et al., 2014. Zhao, R., Ouyang, W., and Wang, X. (2014). Learning mid-level filters for person re-identification. In *CVPR*.

Zheng et al., 2016a. Zheng, L., Bie, Z., Sun, Y., Wang, J., Su, C., Wang, S., and Tian, Q. (2016a). Mars: A video benchmark for large-scale person re-identification. In *ECCV*.

Zheng et al., 2017a. Zheng, L., Huang, Y., Lu, H., and Yang, Y. (2017a). Pose invariant embedding for deep person re-identification. *arXiv:1701.07732*.

Zheng et al., 2015. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., and Tian, Q. (2015). Scalable person re-identification: A benchmark. In *ICCV*.

Zheng et al., 2016b. Zheng, L., Yang, Y., and Hauptmann, A. G. (2016b). Person re-identification: Past, present and future. *arXiv:1610.02984*.

Zheng et al., 2016c. Zheng, L., Zhang, H., Sun, S., Chandraker, M., and Tian, Q. (2016c). Person re-identification in the wild. *arXiv:1604.02531*.

Zheng et al., 2016d. Zheng, Z., Zheng, L., and Yang, Y. (2016d). A discriminatively learned cnn embedding for person re-identification. *arXiv:1611.05666*.

Zheng et al., 2017b. Zheng, Z., Zheng, L., and Yang, Y. (2017b). Unlabeled samples generated by gan improve the person re-identification baseline in vitro. *arXiv:1701.07717*.

Zhong et al., 2017. Zhong, Z., Zheng, L., Cao, D., and Li, S. (2017). Re-ranking person re-identification with k-reciprocal encoding. In *CVPR*.