

Person Re-identification by saliency Learning

Rui Zhao, *Student Member, IEEE*, Wanli Oyang, *Member, IEEE*, and
Xiaogang Wang, *Member, IEEE*

Abstract—Human eyes can recognize person identities based on small salient regions, i.e. human saliency is distinctive and reliable in pedestrian matching across disjoint camera views. However, such valuable information is often hidden when computing similarities of pedestrian images with existing approaches. Inspired by our user study result of human perception on human saliency, we propose a novel perspective for person re-identification based on learning human saliency and matching saliency distribution. The proposed saliency learning and matching framework consists of four steps: (1) To handle misalignment caused by drastic viewpoint change and pose variations, we apply adjacency constrained patch matching to build dense correspondence between image pairs. (2) We propose two alternative methods, i.e. K-Nearest Neighbors and One-class SVM, to estimate a saliency score for each image patch, through which distinctive features stand out without using identity labels in the training procedure. (3) saliency matching is proposed based on patch matching. Matching patches with inconsistent saliency brings penalty, and images of the same identity are recognized by minimizing the saliency matching cost. (4) Furthermore, saliency matching is tightly integrated with patch matching in a unified structural RankSVM learning framework. The effectiveness of our approach is validated on the VIPeR dataset and the CUHK01 dataset. Our approach outperforms the state-of-the-art person re-identification methods on both datasets.

Index Terms—Person re-identification, human saliency, patch matching, video surveillance.

1 INTRODUCTION

Person re-identification [5], [15], [51] is to match pedestrians observed from non-overlapping camera views based on image appearance. It has important applications in video surveillance such as human retrieval, human tracking, and activity analysis. It saves a lot of human efforts on exhaustively searching for a person from large amounts of images and videos. Nevertheless, person re-identification is a very challenging task. A person observed in different camera views undergoes significant variations on viewpoints, poses, and illumination, which make intra-personal variations even larger than inter-personal variations. Image blurring, background clutters and occlusions also cause additional difficulties.

Variations of viewpoints and poses commonly exist in person re-identification, and cause misalignment between images. In Figure 1, the lower right region of ($p1a$) is a red bag, while a leg appears in this region in ($p1b$); the central region of ($p3a$) is an arm, while it becomes a backpack in ($p3b$). Most existing methods [12], [34], [45], [48], [58] match pedestrian images by first computing the difference of feature vectors and then the similarities based on such difference vectors, which is problematic due to the spatial misalignment. In our work, patch matching is employed to handle misalignment, and it is integrated with saliency matching to improve the discriminative power and robustness to spatial variation.

Salient regions in pedestrian images provide valuable information in identification. However, if they are small

in size, saliency information is often overwhelmed by other features when computing similarities of images. In this paper, *saliency* means regions with attributes that 1) make a person *distinctive* from their candidates, and 2) are *reliable* in finding the same person across camera views. In many cases, humans can easily recognize matched pedestrian pairs because they have distinct features. For example, in Figure 1, person $p1$ takes a red bag, $p2$ dresses bright white skirt, $p3$ takes a blue bag, and $p4$ carries a red folder in arm. These features are discriminative in distinguishing one person from others. Intuitively, if a body part is salient in one camera view, it usually remains salient in another camera view. Therefore, saliency also has view invariance.

Salient regions are not limited to body parts (such as clothes and trousers), but also include accessories (such as baggages, folders and umbrellas as shown in Figure 1), which are often considered as outliers and removed in existing approaches. Our computation of saliency is based on the comparison with images from a large scale reference dataset rather than a small group of persons. Therefore, it is quite stable in most circumstances.

We observe that images of the same person captured from different camera views have some invariance property on their spatial distributions of saliency, like pair ($a1, a2$) in Figure 2. Since the person in image ($a1$) shows saliency in her dress while others ($a3$)-($a6$) are salient in blouses, they can be well distinguished simply from the spatial distributions of saliency. Therefore, not only the visual features from salient regions are discriminative, the spatial distributions of human saliency also provide useful information in person re-identification. Such information can be encoded into patch matching. If two patches from two images of the same person are

• R. Zhao, W. Ouyang, and X. Wang are with the Department of Electronic Engineering, the Chinese University of Hong Kong, Hong Kong.
E-mail: {rzhao, wlouyang, xgwang}@ee.cuhk.edu.hk

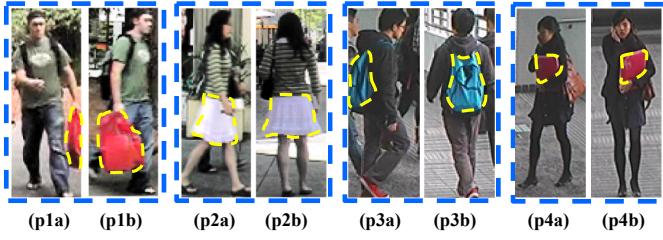


Fig. 1. Salient region could be a body part or a carrying accessory. Some salient regions of pedestrians are highlighted with yellow dashed boundaries.

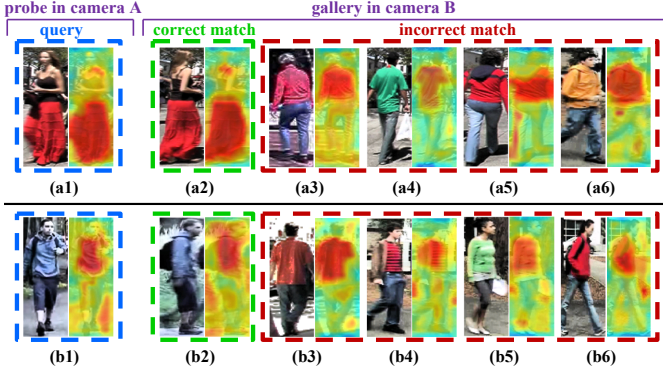


Fig. 2. Illustration of saliency matching with examples. saliency map of each pedestrian image is shown. **Best viewed in color.**

matched, they are expected to have similar saliency values; otherwise such matching brings penalty on saliency matching. In the second row in Figure 2, the query image (b1) shows a similar saliency distribution as those of gallery images. In this case, visual similarity needs to be considered. This motivates us to relate saliency matching penalty to the visual similarity of two matched patches.

2 OUR APPROACH

Although saliency plays an important role in person re-identification, it has not been well explored in literature. In this paper, a novel framework of human saliency learning and matching is proposed for person re-identification. Our major contributions can be summarized from the following aspects.

- 1) We propose a way of estimating saliency based on human perception through user study. It is estimated from the number of trials that a human subject recognizes a query image from a candidate pool only based on a local region. It shows that most pedestrian images can be matched by humans from local salient regions without looking at whole images. The saliency estimated from user study is compared with the result of our saliency computation model. Compared with general image saliency detection methods [6], [14], our proposed saliency computation has much stronger correlation with human perception in person re-identification.

- 2) A computation model is proposed to estimate the probabilistic saliency map. Different from general image saliency detection methods, it is specially designed for person re-identification, and has the following properties. 1) It is robust to changes of viewpoints, poses and articulation. 2) Distinct patches are considered as salient only when they are matched and distinct in both camera views. 3) Human saliency itself is a useful descriptor for pedestrian matching. For example, a person only with salient upper body and a person only with salient lower body must be different identities.
- 3) We formulate person re-identification as a saliency matching problem. Dense correspondences between patches are established by patch matching based on visual similarity, and matching patches with inconsistent saliency brings cost. Images of the same person are recognized by minimizing the saliency matching cost, which depends on both locations and visual similarity of matched patches.
- 4) saliency matching and patch matching are tightly integrated into a unified structural RankSVM framework. Structural RankSVM has good training efficiency given a large number of rank constraints in person re-identification. Our approach transforms the original high-dimensional visual feature space to a 80 times lower dimensional saliency feature space to further improve training efficiency and also avoid overfitting.

3 RELATED WORKS

Existing works on person re-identification mainly focus on two aspects: 1) *features and representations*, and 2) *distance metric*. A review can be found in [15].

3.1 Features and Representations

A lot of research efforts [4], [10], [11], [13], [41]–[43], [52], [54], [57], [59] have been devoted to exploiting discriminative features in person re-identification. Wang *et al.* [52] proposed shape and appearance context to model the spatial distributions of appearance relative to body parts in order to extract discriminative features robust to misalignment. Farenzena *et al.* [13] proposed the Symmetry-Driven Accumulation of Local Features (SDALF) by exploiting the symmetry property in pedestrian images to handle view variation. Bak *et al.* [4], Xu *et al.* [54] and Cheng *et al.* [10], [11] applied human part models and pictorial structures to cope with pose variations by establishing the spatial correspondence. Ma *et al.* [41]–[43] developed the BiCov descriptor based on the Gabor filters and the covariance descriptor to handle illumination change and background variation. Zheng *et al.* [57], [59] used the contextual visual cues from surrounding people to enrich human signatures. Information on salient regions exploited in our work can be integrated with many of these feature designs by putting more weights on features from salient regions.

Features vary in their usefulness in person matching, and some works have been done on feature selection and importance learning. Gray *et al.* [17] used AdaBoost to select features. Schwartz [49] assigned weights to features with Partial Least Squares (PLS). Liu *et al.* [37] developed an unsupervised approach to learn bottom-up feature importance, and adaptively weight features. Instead of globally weighting features across all the pedestrian images, our approach adaptively weights features based on individual person pairs to be matched, since different persons have different salient regions.

Visual features suffer from a range of variations across camera views. Feature transforms are learned to improve the invariance to cross-view transforms. Prosser *et al.* [47] learned the Cumulative Brightness Transfer Function to handle color transforms. Avraham *et al.* [2], [3] learned both implicit and explicit transforms of visual features. Li and Wang [33] learned a mixture of cross-view transforms and projected features into a common space for alignment. Rather than learning feature transforms for specific camera view settings, our approach flexibly handle the cross-view variations by performing a constrained patch matching technique, which can be generalized to any disjoint camera-view transition.

Some works explored higher level features [28]–[30], [39], [50], [56] to assist person re-identification. Vaquero *et al.* [50] first introduced mid-level facial attributes in human recognition. Layne *et al.* [28], [29] proposed 15 human attributes for person re-identification. Song *et al.* [39] used human attributes to match persons with Bayesian decision. saliency distribution can also be considered as one kind of high-level features.

3.2 Rank and Metric Learning

Given a query image, an image of the same person is expected to have a high rank on the candidate list after matching. Prosser *et al.* [48] formulated person re-identification as a ranking problem, and learned global feature weights with RankSVM. Wu *et al.* [53] introduced rank-loss optimization to improve accuracy in re-identification. Loy *et al.* [40] exploited unlabeled gallery data to propagate labels to query instances with a manifold ranking model. Liu *et al.* [38] presented a man-in-loop method to allow users to quickly refine ranking result. In this paper, we employ structural RankSVM [24], which considers ranking difference.

Many research works [12], [19], [20], [26], [34], [35], [37], [45], [46], [58] focused on optimizing distance metrics for matching persons. Zheng *et al.* [58] learned the metric by maximizing the likelihood of true matches to have a smaller distance than that of a wrongly matched pair. Dikmen *et al.* [12] proposed to learn a Mahalanobis distance that is optimal for k-nearest neighbor classification by using a maximum margin formulation. Mignon and Jurie [45] learned a joint projection for dimension reduction, satisfying distance constraints added by image pairs. Li *et al.* [35] proposed to learn a decision function

for matching, which jointly models a distance metric and a locally adaptive thresholding rule. Pedagadi *et al.* [46] employed Local Fisher Discriminant Analysis to learn a distance metric. Above learned metrics are based on subtraction of misaligned feature vectors, which causes significant information loss and errors. Our approach handles feature misalignment through patch matching.

3.3 Human saliency vs. General Image saliency

General image saliency has been well studied [14], [21], [22], [25], [32]. In the context of person re-identification, human saliency is different than general image saliency in the way of drawing visual attentions. With the aim to improve the performance of re-identification, human saliency is considered as visual patterns of distinguishing a person from others, while general saliency draws visual attention within a single image to capture salient foreground objects from background.

4 METHOD OVERVIEW

The diagram of the proposed saliency learning and matching framework is shown in Figure 3. Section 5 conducts a user study to estimate human saliency based on human perception in the person re-identification task. We investigate the discriminative power of different body regions in identifying a target person from a gallery set. The saliency of each local region of a query image is quantitatively estimated by measuring the averaged number of trails that human labelers find the target person only based on that region of the query image. An illustration is shown in Figure 3 (a). The red and green bounding boxes indicate incorrect and correct targets chosen by the labeler from the gallery. The red skirt has higher saliency and causes fewer failure trails compared with the arm. Our result shows that subjects can recognize a query person only based on a small salient part without looking at the whole image. Salient regions vary on different persons.

An unsupervised approach for saliency learning is proposed in Section 6 and illustrated in Figure 3 (b). With constrained patch matching, each patch finds its matched neighbors from a reference set of training images. K-Nearest Neighbor and One-Class SVM models are employed to learn human saliency. Our experimental results show both qualitative and quantitative evaluation of the correlation between the learned saliency and human perception. With obtained human saliency, matching image pairs can be performed in unsupervised and supervised ways as described in Section 6. For the unsupervised manner, saliency is used to bi-directionally weight patch matching similarity and penalize inconsistency of saliency distribution across camera views, as shown by the blue lines in Figure 3 (c). For the supervised manner, person matching is formulated as a saliency matching problem, which considers four types of saliency matching cases, as shown in the table in Figure 3 (c). Matching cost is a linear function of patch

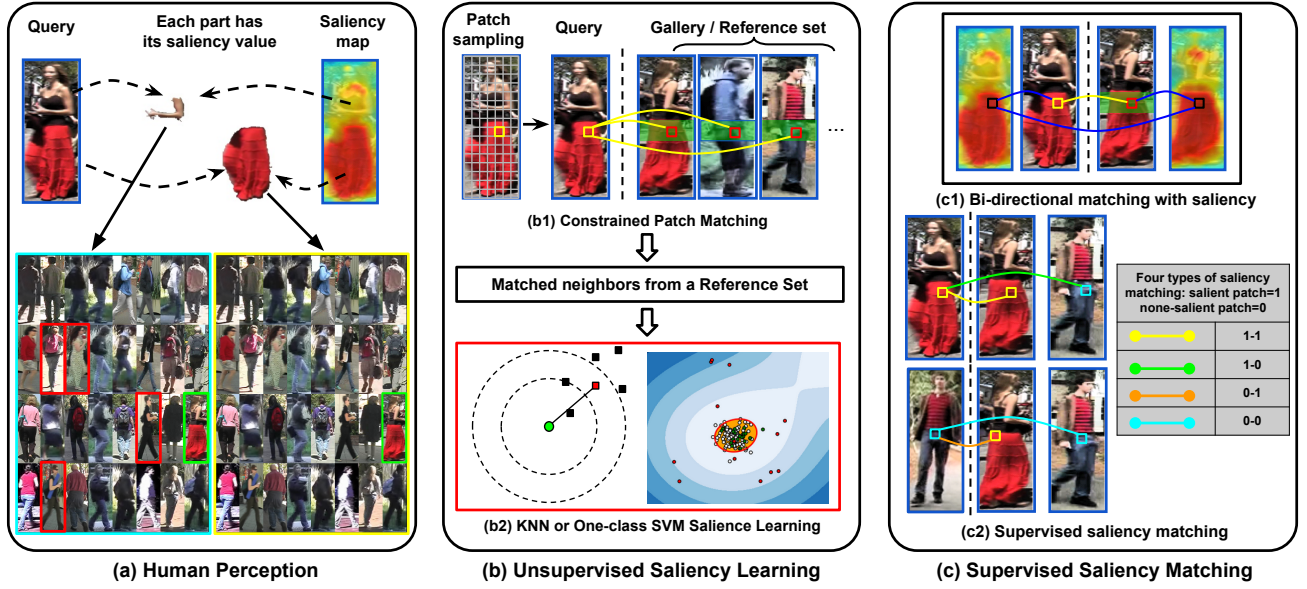


Fig. 3. Diagram of our novel framework of human saliency learning and matching for person re-identification.



Fig. 4. Salient and non-salient body parts in person re-identification.

matching similarities, which is learned with Structural RankSVM. The learned saliency matching function is used to measure similarities between images.

5 SALIENCY FROM HUMAN PERCEPTION

We define human saliency in the context of person re-identification and estimate it by user study. The design of user study considers the following aspects.

- Salient body parts are those possessing unique appearance compared with a reference set.
- Human body, including carrying accessories, can be decomposed into parts with different saliency values.
- If a body part helps subjects to quickly identify a person from other candidates across camera views, it is considered as salient.
- The salient values of different parts are estimated independently to simplify analysis. Higher order saliency from combinations of body parts could be studied in the future work.

As an example in Figure 4, the yellow bag is a carrying accessory, which can be easily identified across camera

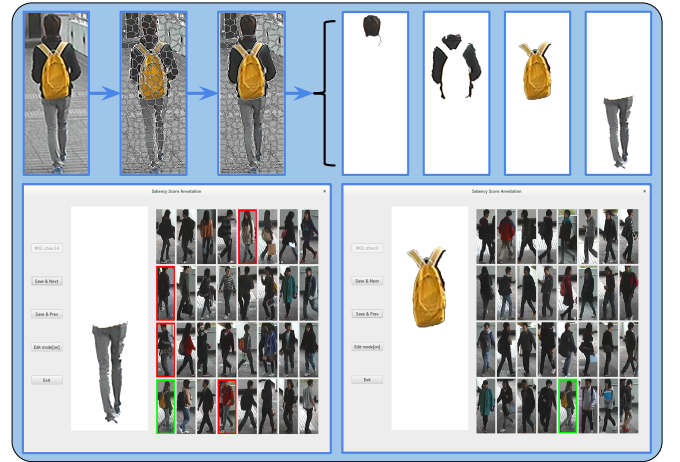


Fig. 5. Flow chart of human saliency annotation. The first row illustrates the procedure that an image is segmented into semantic body parts. The second row shows the interface of annotating human saliency.

views, while the black coat appears on many persons, and is hard to be used as a cue to recognize identity. Thus, the yellow bag has a higher saliency value.

5.1 Human Annotation Scheme

Given an image, we apply superpixel segmentation [1], and then manually merge superpixels that are coherent in appearance. As an example shown in the first row of Figure 5, superpixels on the yellow bag are merged into a part. Superpixels with different semantic meanings are not merged. For example, even if the hair and jacket share similar appearance, they are annotated as two parts. Only foreground superpixels are annotated.

A segmented body part is randomly selected and presented to a labeler for annotation. Each part is annotated



Fig. 6. **Examples of saliency annotation.** Each body part is annotated with a saliency value. saliency map is overlaid on the gray-level image (right). The original color image is on the left.

for multiple times by different labelers. Their annotations are combined into a saliency value. In Figure 5, a body part from a query image is revealed (on the left) at its original spatial location in the image while other parts are masked, and a list of 32 images randomly sampled from the gallery set are also shown (on the right) to the labeler. The true target (observed in a different camera view from the query image) is among the sampled images, but the order is randomly shuffled. The labeler is asked to select the most likely image from the list based on visual perception. The labeler is allowed to select for multiple times until the correct match is found. In the second row of Figure 5, the red bounding boxes indicate wrong selection and the green one indicates the correct match found in the end. A part is considered as salient if labelers try fewer times to found the target.

Denote the i -th revealed part by p_i , and the number of trails of the j -th user on this part by $n_{p_i}^j$. Then the saliency value of the revealed part is estimated as

$$\text{score}(p_i) = \exp\left(-\frac{m_{p_i}^2}{\sigma_{avg}^2}\right) \exp\left(-\frac{s_{p_i}^2}{\sigma_{std}^2}\right). \quad (1)$$

m_{p_i} and s_{p_i} are the average and standard deviation of $n_{p_i}^j$ over all the labelers. σ_{avg} and σ_{std} are bandwidth parameters.

Annotation is conducted on 524 body parts of 100 images from camera view A of the VIPeR dataset [16]. Some examples of the annotated saliency maps are shown in Figure 6. In order to investigate whether salient regions exist in pedestrian images, Figure 7 (a) shows the histogram on the numbers of trails used to find the targets only based on the most salient parts on query images. It shows that more than half of the pedestrians can be recognized, if the labelers only observe the most salient part of a query image. As comparison, Figure 7 (b) plots the histogram on the numbers of trails for all the parts. It shows that most other body parts are not salient enough. The correlation between the annotated saliency and that obtained with the proposed computation model will be validated in experiments in Section 8.

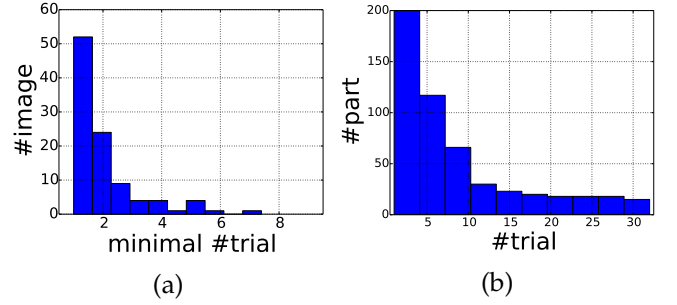


Fig. 7. **Statistics on saliency annotation.** (a) Histogram on the numbers of trails used to find the targets only based on the most salient parts on query images. (b) Histogram on the numbers of trails for all the parts.

6 HUMAN SALIENCY LEARNING

We propose to automatically learn human saliency in an unsupervised manner. Dense correspondence between images is first built with patch matching, and two alternative approaches (K-nearest neighbor and One-Class SVM) are proposed to estimate human saliency without using identity labels or human annotated saliency.

6.1 Feature Extraction

Each image is densely divided into a grid of overlapping local patches, and each patch is represented by a feature vector concatenating color histograms and SIFT features computed around its local region.

Dense Color Histogram. A color histogram in LAB color space is extracted from each patch. LAB color histograms are computed on multiple downsampled scales and L2 normalized.

Dense SIFT. We divide each patch into 4×4 cells, quantize the orientations of local gradients into 8 bins, and obtain a $4 \times 4 \times 8 = 128$ dimensional SIFT feature vector, which is also L2 normalized.

In our experiment, scales of pedestrian images range from 128×48 to 160×60 . Patches of size 10×10 pixels are sampled on a dense grid with a step size 4. 32-bin color histograms are computed in each LAB channels, and in each channel, 3 levels of downsampling are used with scaling factors 0.5, 0.75 and 1. SIFT features are also extracted in 3 color channels and thus produces a 128×3 feature vector for each patch. In a summary, each patch is finally represented with a discriminative descriptor vector of length $32 \times 3 \times 3 + 128 \times 3 = 672$. We denote the combined Color-SIFT feature vector as *DenseFeats*.

6.2 Dense Correspondence

To deal with misalignment, we build dense correspondence between images by adjacency constrained search. DenseFeats features of a pedestrian image is represented as $X^{A,u} = \{\mathbf{x}_{m,n}^{A,u} \mid m = 1 \dots, M, n = 1 \dots, N\}$, where (A, u) denotes the u -th image in camera A , (m, n) denotes the patch centered at the m -th row and the n -th column of this image, and $\mathbf{x}_{m,n}^{A,u}$ is the dense Color-SIFT feature



Fig. 8. Illustration of adjacency constrained search. Green region represents the adjacency constrained search set of patch in yellow box. Patch in red box is the target match.

vector of the patch. A natural baseline is to compute image similarity with concatenated patch features,

$$\text{sim}_{\text{DenseFeats}}(X^{A,u}, X^{B,v}) = \sum_{i=1, \dots, M} \sum_{j=1, \dots, N} s(\mathbf{x}_{i,j}^{A,u}, \mathbf{x}_{i,j}^{B,v}), \quad (2)$$

where

$$s(\mathbf{x}_{i,j}^{A,u}, \mathbf{x}_{i,j}^{B,v}) = \exp\left(-\frac{d(\mathbf{x}_{i,j}^{A,u}, \mathbf{x}_{i,j}^{B,v})^2}{2\sigma^2}\right), \quad (3)$$

is the similarity between two patch features, $d(\cdot)$ is the Euclidean distance, and σ is a bandwidth parameter.

Adjacency Searching. $\text{sim}_{\text{DenseFeats}}$ does not consider misalignment. We propose adjacency constrained searching to allow flexible matching among patches in image pairs. When the patches are matched with those from another image, patches in the same row have the same search set, denoted as

$$\mathcal{S}(\mathbf{x}_{m,n}^{A,u}, X^{B,v}) = \{\mathbf{x}_{i,j}^{B,v} \mid i = m, j = 1, \dots, N\}. \quad (4)$$

$\mathcal{S}(\mathbf{x}_{m,n}^{A,u}, X^{B,v})$ restricts the search set in $X^{B,v}$ within the m -th row. However, bounding boxes produced by a human detector are not always well aligned, and also uncontrolled human pose variations exist. We relax the horizontal constraint to have a larger search range:

$$\hat{\mathcal{S}}(\mathbf{x}_{m,n}^{A,u}, X^{B,v}) = \{\mathbf{x}_{i,j}^{B,v} \mid i \in \mathcal{N}(m), j = 1, \dots, N\}, \quad (5)$$

where $\mathcal{N}(m) = \{\max(0, m-l), \dots, m, \dots, \min(m+l, M)\}$. l defines the size of the relaxed adjacent vertical space. Less relaxed search space cannot well tolerate the spatial variation while more relaxed search space increases the chance of matching different body parts. $l = 2$ is chosen in our setting.

We perform the nearest neighbor search for each $\mathbf{x}_{m,n}^{A,p}$ in its search set $\hat{\mathcal{S}}(\mathbf{x}_{m,n}^{A,p}, X^{B,q})$. For each patch $\mathbf{x}_{m,n}^{A,p}$, a nearest neighbor is sought from its search set in every image within a reference set. The adjacency constrained search process is illustrated in Figure 8, and some visually similar patches returned by the discriminative adjacency constrained search are shown in Figure 9.

6.3 Unsupervised saliency Learning

6.3.1 K-Nearest Neighbor (KNN) saliency

Byers *et al.* [7] found the KNN distances can be used for clutter removal. Since human saliency detection shares a similar goal as abnormality detection, KNN should also

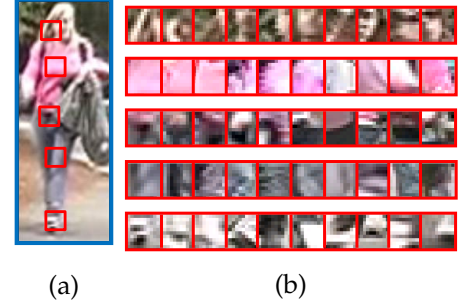


Fig. 9. **Examples of adjacency search.** (a) A test image from the VIPeR dataset. Local patches are densely sampled, and five exemplar patches on different body parts are shown in red boxes. (b) One nearest neighbor from each reference image is returned by adjacency search for each patch on the left, and then N nearest neighbors from N reference images are sorted. The top ten nearest neighbor patches are shown. Note that the ten nearest neighbors are from ten different images.

be viable in finding human saliency. By searching for the K -nearest neighbors of a test patch in the set of matched patches obtained with dense correspondence, KNN is adapted to the re-identification problem. saliency score of the test patch is computed with the KNN distance.

We denote the number of images in the reference set by N_r . After building dense correspondences between a test image and reference images, a nearest neighbor (NN) set of size N_r is obtained for every patch $\mathbf{x}_{m,n}^{A,u}$,

$$X_{NN}(\mathbf{x}_{m,n}^{A,u}) = \{\mathbf{x} \mid \underset{\mathbf{x}_{i,j}^{B,v}}{\operatorname{argmin}} d(\mathbf{x}_{m,n}^{A,u}, \mathbf{x}_{i,j}^{B,v}), \mathbf{x}_{i,j}^{B,v} \in \hat{\mathcal{S}}(\mathbf{x}_{m,n}^{A,u}, X^{B,v}), v = 1, 2, \dots, N_r\} \quad (6)$$

The KNN distances between $\mathbf{x}_{m,n}^{A,u}$ and its nearest neighbors in $X_{NN}(\mathbf{x}_{m,n}^{A,u})$ are used as the saliency score:

$$\text{score}_{knn}(\mathbf{x}_{m,n}^{A,u}) = d_k(X_{NN}(\mathbf{x}_{m,n}^{A,u})), \quad (7)$$

where d_k denotes the distance of the k -th nearest neighbor. Salient patches only find a limited number ($k = \alpha_k N_r$) of visually similar neighbors, as shown in Figure 10, and then $\text{score}_{knn}(\mathbf{x}_{m,n}^{A,p})$ is expected to be large. $0 < \alpha_k < 1$ is a proportion parameter reflecting our expectation on the statistical distribution of salient patches.

Choosing k . The goal of saliency detection for person re-identification is to identify parts with unique appearance. We set $\alpha_k = 0.5$ with an empirical assumption that a patch is considered to have unique appearance such that more than half of the people in the reference set do not share similar patches with it. N_r reference images are randomly sampled from training set in our experiments. Enlarging the reference dataset will not deteriorate saliency detection, because saliency is defined in the statistical sense. It is robust as long as the distribution of the reference dataset well reflects the test scenario.

6.3.2 One-Class SVM saliency

One-class SVM [18] has been widely used for outlier detection. The basic idea is to use a hypersphere to



Fig. 10. **Illustration of salient patch distribution.** Salient patches are distributed far away from other patches.

describe data in the feature space and put most of the data into the hypersphere. It is formulated as an objective function:

$$\min_{R \in \mathbb{R}, \xi \in \mathbb{R}^l, c \in F} R^2 + \frac{1}{vl} \sum_i \xi_i, \quad (8)$$

$$s.t. \|\Phi(\mathbf{x}_i) - c\|^2 \leq R^2 + \xi_i, \quad \forall i \in \{1, \dots, l\} : \xi_i \geq 0,$$

where $\Phi(\mathbf{x}_i)$ is the multi-dimensional feature vector of i -th training sample, l is the number of training samples, R and c are the radius and center of the hypersphere, and $v \in [0, 1]$ is a trade-off parameter. The goal is to keep the hypersphere as small as possible and include most of the training data. It can be solved in a dual form by QP optimization [9]. The decision function is:

$$f(\mathbf{x}) = R^2 - \|\Phi(\mathbf{x}) - c\|^2, \quad (9)$$

$$\|\Phi(\mathbf{x}) - c\|^2 = k(\mathbf{x}, \mathbf{x}) - 2 \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) + \sum_{i,j} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j),$$

where α_i and α_j are the parameters for each constraint in the dual problem. We use the radius basis function (RBF) $K(\mathbf{x}, \mathbf{y}) = \exp\{-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2\}$ as kernel to deal with high-dimensional, non-linear, and multi-mode distributions. As shown in [9], the decision function of kernel One-class SVM can well capture the density and modality of the feature distribution. As an alternative to the KNN saliency algorithm (Section 6.3.1) without requiring the choice of K , saliency score is defined in terms of kernel One-class SVM decision function:

$$\begin{aligned} \text{score}_{ocsvm}(\mathbf{x}_{m,n}^{A,u}) &= d(\mathbf{x}_{m,n}^{A,u}, \mathbf{x}^*), \\ \mathbf{x}^* &= \underset{\mathbf{x} \in X_{NN}(\mathbf{x}_{m,n}^{A,u})}{\operatorname{argmax}} f(\mathbf{x}). \end{aligned} \quad (10)$$

Our experiments show very similar results in person re-identification with both saliency detection methods. score_{ocsvm} performs slightly better than score_{knn} in some circumstances. The probability of $\mathbf{x}_{m,n}^{A,u}$ being a salient patch is

$$P(l_{m,n}^{A,u} = 1 | \mathbf{x}_{m,n}^{A,u}) = 1 - \exp(-\text{score}_{opt}(\mathbf{x}_{m,n}^{A,u})^2 / \sigma_0^2), \quad (11)$$

where $opt \in \{knn, ocsvm\}$. The human saliency learning is summarized in Algorithm 1.

Algorithm 1 Human saliency learning.

Input: image $X^{A,u}$ and a reference image set $\mathcal{R} = \{X^{B,v}, v = 1, \dots, N_r\}$

Output: saliency probability map $P(l_{m,n}^{A,u} = 1 | \mathbf{x}_{m,n}^{A,u})$

- 1: **for** each patch $\mathbf{x}_{m,n}^{A,u} \in X$ **do**
 - 2: compute $X_{NN}(\mathbf{x}_{m,n}^{A,u})$ with Eq. (6)
 - 3: compute $\text{score}_{opt}(\mathbf{x}_{m,n}^{A,u})$, $opt \in \{knn, ocsvm\}$ with Eq. (7) or Eq. (10)
 - 4: compute $P(l_{m,n}^{A,u} = 1 | \mathbf{x}_{m,n}^{A,u})$ with Eq. (11)
 - 5: **end for**
-

7 SALIENCY MATCHING

One of our main contributions is to match human images based on saliency probability map. It is based on our observation that person in different camera views shows consistence in saliency probability maps, as shown in Figure 2. Since matching is applied to arbitrary image pairs, we omit the image index in notation for concise clarity, *i.e.* change $X^{A,u}$ to X^A , $X^{B,v}$ to X^B , $\mathbf{x}_{m,n}^{A,u}$ to $\mathbf{x}_{p_i}^A$ and $\mathbf{x}_{i,j}^{B,v}$ to $\mathbf{x}_{p'_i}^B$. p_i is the patch index in image X^A and p'_i is the corresponding matched patch index in image X^B produced by dense correspondence. We denote the dense correspondence between X^A and X^B as $P = \{(p_i, p'_i)\}_{i=1, \dots, MN}$.

7.1 Bi-directional Weighted Matching

We first denote the method of only using patch matching without saliency information as *PatMatch*, and the image similarity is expressed as

$$\text{sim}_{PatMatch}(X^A, X^B) = \sum_{(p_i, p'_i) \in P} s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B). \quad (12)$$

$s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B)$ is the visual similarity between patches. Searching for the best matched image in the gallery is formulated as finding the maximal similarity score.

$$v^* = \underset{v}{\operatorname{argmax}} \text{sim}(X^{A,u}, X^{B,v}) \quad (13)$$

A bi-directional weighted matching is designed to incorporate saliency information. We denote this method as saliency guided dense correspondence (*SDC*), as illustrated in Figure 3(c1), and the similarity between two images is computed as

$$\text{sim}_{SDCopt} = \sum_{(p_i, p'_i) \in P} \frac{\text{score}_{opt}(\mathbf{x}_{p_i}^A) \cdot s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) \cdot \text{score}_{opt}(\mathbf{x}_{p'_i}^B)}{\alpha_{sdc} + |\text{score}_{opt}(\mathbf{x}_{p_i}^A) - \text{score}_{opt}(\mathbf{x}_{p'_i}^B)|}, \quad (14)$$

where α_{sdc} is a parameter representing a base penalty. Intuitively, large saliency scores in both matched patches are expected to enhance the similarity score of matched patches. In another aspect, images of the same person would be more likely to have similar saliency distributions than those of different persons, so the difference in saliency score can be used as a penalty to the similarity score. We set $\alpha_{sdc} = 1$ in experiments.

7.2 Unified saliency Matching

To incorporate saliency into matching, we introduce $L^A = \{l_{p_i}^A \mid l_{p_i}^A \in \{0, 1\}\}$ and $L^B = \{l_{p'_i}^B \mid l_{p'_i}^B \in \{0, 1\}\}$ as saliency labels for all the patches in image X^A and X^B respectively. If all the saliency labels are known, we can perform person matching by computing the saliency matching score as follows:

$$f_z(X^A, X^B, L^A, L^B; P, Z) = \sum_{(p_i, p'_i) \in P} \left\{ z_{p_i,1} l_{p_i}^A l_{p'_i}^B + z_{p_i,2} l_{p_i}^A (1 - l_{p'_i}^B) + z_{p_i,3} (1 - l_{p_i}^A) l_{p'_i}^B + z_{p_i,4} (1 - l_{p_i}^A) (1 - l_{p'_i}^B) \right\}, \quad (15)$$

where $Z = \{z_{p_i,k} \mid i=1, \dots, MN, k=1, 2, 3, 4\}$ are the matching scores for four different saliency matching results at one local patch. $z_{p_i,k}$ is not a constant for all the patches. Instead, it depends on the spatial location p_i . For example, the score of matching patches on the background should be different than those on legs. $z_{p_i,k}$ also depends on the visual similarity between patches $\mathbf{x}_{p_i}^A$ and patch $\mathbf{x}_{p'_i}^B$. Instead of directly using the Euclidean distance $d(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B)$, we convert it to similarity to reduce the side effect in summation of very large distances in incorrect matching, caused by misalignment, occlusion, or background clutters.

Therefore, we define the matching score $z_{p_i,k}$ as a linear function of the similarity as follows,

$$z_{p_i,k} = \alpha_{p_i,k} \cdot s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) + \beta_{p_i,k}. \quad (16)$$

$\alpha_{p_i,k}$ and $\beta_{p_i,k}$ are weighting parameters. Thus Eq. (15) considers both saliency matching and visual similarity.

Since the saliency labels $l_{p_i}^A$ and $l_{p'_i}^B$ in Eq. (15) are hidden variables, they can be marginalized by computing the expectation of the saliency matching score as

$$\begin{aligned} f^*(X^A, X^B; P, Z) &= \sum_{L^A, L^B} f_z(X^A, X^B, L^A, L^B; P, Z) p(L^A, L^B \mid X^A, X^B) \\ &= \sum_{(p_i, p'_i) \in P} \sum_{k=1}^4 \left[\alpha_{p_i,k} \cdot s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) + \beta_{p_i,k} \right] c_{p_i,k}(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B), \end{aligned} \quad (17)$$

where $c_{p_i,k}(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B)$ is the probabilistic saliency matching cost depending on saliency probabilities $P(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A)$ and $P(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B)$ given in Eq. (11),

$$c_{p_i,k}(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) = \begin{cases} p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B), & k = 1, \\ p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B), & k = 2, \\ p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B), & k = 3, \\ p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B), & k = 4. \end{cases} \quad (18)$$

To better formulate this learning problem, we extract out

all the weighting parameters in Eq. (17) as \mathbf{w} , and have

$$\begin{aligned} f^*(X^A, X^B; P, Z) &= \mathbf{w}^T \Phi(X^A, X^B; P) \\ &= \sum_{(p_i, p'_i) \in P} \mathbf{w}_{p_i}^T \phi(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B), \end{aligned} \quad (19)$$

where

$$\begin{aligned} \Phi(X^A, X^B; P) &= [\phi(\mathbf{x}_{p_1}^A, \mathbf{x}_{p'_1}^B)^T, \dots, \phi(\mathbf{x}_{p_{MN}}^A, \mathbf{x}_{p'_{MN}}^B)^T]^T, \\ \mathbf{w} &= [\mathbf{w}_{p_1}, \dots, \mathbf{w}_{p_{MN}}]^T, \\ \mathbf{w}_{p_i} &= [\{\alpha_{p_i,k}\}_{k=1,2,3,4}, \{\beta_{p_i,k}\}_{k=1,2,3,4}]. \end{aligned} \quad (20)$$

$\Phi(X^A, X^B; P)$ is the feature map describing the matching between X^A and X^B . For each patch p_i , the matching feature $\phi(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B)$ is an eight dimensional vector:

$$\phi(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) = \begin{bmatrix} s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) \cdot p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B) \\ s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) \cdot p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B) \\ s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) \cdot p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B) \\ s(\mathbf{x}_{p_i}^A, \mathbf{x}_{p'_i}^B) \cdot p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B) \\ p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B) \\ p(l_{p_i}^A = 1 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B) \\ p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 1 \mid \mathbf{x}_{p'_i}^B) \\ p(l_{p_i}^A = 0 \mid \mathbf{x}_{p_i}^A) \cdot p(l_{p'_i}^B = 0 \mid \mathbf{x}_{p'_i}^B) \end{bmatrix}. \quad (21)$$

As shown in Eq. (21), the pairwise feature map $\Phi(X^A, X^B; P)$ combines the saliency probability map with appearance matching similarities. For each query image X^A , the images in the gallery are ranked according to the expectations of saliency matching scores in Eq. (17). There are three advantages of matching with human saliency : (1) the human saliency probability distribution is more invariant than other features in different camera views; (2) because the saliency probability map is built based on dense correspondence, it inherits the property of tolerating spatial variation; and (3) it can be weighted by visual similarity to improve the performance of person re-identification. We will present the details in next section by formulating the person re-identification problem with $\Phi(X^A, X^B; P)$ in the structural RankSVM framework.

7.3 Ranking by Partial Order

We cast person re-identification as a ranking problem for supervised training. The ranking problem will be solved by finding an optimal partial order, mathematically defined in Eq. (22)(23)(26). Given a dataset of pedestrian images, $\mathcal{D}^A = \{X^{A,u}, id^{A,u}\}_{u=1}^U$ from camera view A and $\mathcal{D}^B = \{X^{B,v}, id^{B,v}\}_{v=1}^V$ from camera view B , where $X^{A,u}$ is the u -th image, $id^{A,u}$ is its identity label, and U is the total number of images in \mathcal{D}^A . Similar notations apply for variables of camera view B . Each image $X^{A,u}$ has its relevant images (same identity) and irrelevant images (different identities) in dataset \mathcal{D}^B . Our goal is to learn the weight parameters \mathbf{w} that order relevant gallery images before irrelevant ones. For the image

$X^{A,u}$, we rank the relevant images before irrelevant ones, but no information of the orders within relevant images or irrelevant ones is provided. The partial order $\mathbf{y}^{A,u}$ is denoted as,

$$\mathbf{y}^{A,u} = \{y_{v,v'}^{A,u}\}, \quad y_{v,v'}^{A,u} = \begin{cases} +1 & X^{B,v} \prec X^{B,v'} \\ -1 & X^{B,v} \succ X^{B,v'} \end{cases}, \quad (22)$$

where $X^{B,v} \prec X^{B,v'}$ ($X^{B,v} \succ X^{B,v'}$) represents that $X^{B,v}$ is ranked before (after) $X^{B,v'}$ in partial order $\mathbf{y}^{A,u}$.

The partial order feature [23], [44] is appropriate for our goal and can encode the difference between relevant pairs and irrelevant pairs with only partial orders. The partial order feature for image $X^{A,u}$ is formulated as,

$$\Psi_{po}(X^{A,u}, \mathbf{y}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V) = \sum_{X^{B,v} \in S_{X^{A,u}}^+} \sum_{X^{B,v'} \in S_{X^{A,u}}^-} \frac{y_{v,v'}^{A,u} \Phi(X^{A,u}, X^{B,v}; P^{u,v}) - \Phi(X^{A,u}, X^{B,v'}; P^{u,v'})}{|S_{X^{A,u}}^+| \cdot |S_{X^{A,u}}^-|}, \quad (23)$$

$$S_{X^{A,u}}^+ = \{X^{B,v} \mid id^{B,v} = id^{A,u}\}, \quad (24)$$

$$S_{X^{A,u}}^- = \{X^{B,v} \mid id^{B,v} \neq id^{A,u}\}, \quad (25)$$

where $\{P^{u,v}\}_{v=1}^V$ are the dense correspondences between image $X^{A,u}$ and every gallery image $X^{B,v}$, $S_{X^{A,u}}^+$ is relevant image set of $X^{A,u}$, $S_{X^{A,u}}^-$ is irrelevant image set, $\Phi(X^{A,u}, X^{B,v}; P^{u,v})$ is the feature map defined in Eq. (20), and the difference vector of two feature maps $\Phi(X^{A,u}, X^{B,v}; P^{u,v}) - \Phi(X^{A,u}, X^{B,v'}; P^{u,v'})$ is added if $X^{B,v} \prec X^{B,v'}$ or subtracted otherwise.

A partial order may correspond to multiple rankings. Our task is to find a good ranking satisfying the optimal partial order $\mathbf{y}_*^{A,u}$ that maximizes the following score function,

$$\mathbf{y}_*^{A,u} = \underset{\mathbf{y}^{A,u} \in \mathcal{Y}^{A,u}}{\operatorname{argmax}} \quad \mathbf{w}^T \Psi_{po}(X^{A,u}, \mathbf{y}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V), \quad (26)$$

where $\mathcal{Y}^{A,u}$ is the space consisting of all the possible partial orders. As discussed in [23], [55], good ranking can be obtained by sorting gallery images by $\{\mathbf{w}^T \Phi(X^{A,u}, X^{B,v}; P^{u,v})\}_v$ in a descending order. The remaining problem is how to learn \mathbf{w} . With an optimized \mathbf{w}_* , we denote the unified saliency matching similarity as

$$\operatorname{sim}_{SalMatch_{opt}}(X^A, X^B) = \mathbf{w}_*^T \Phi(X^A, X^B; P), \quad (27)$$

where $opt \in \{knn, ocsvm\}$.

7.4 Structural RankSVM Training

We employ structural SVM to learn the weighting parameters \mathbf{w} . Different than many previous SVM-based approaches [8], [48] doing optimization over pairwise differences, structural SVM optimizes over ranking differences and can incorporate non-linear multivariate loss functions into global optimization in SVM training.

Objective function. Our goal is to learn a linear model and the training is based on n-slack structural SVM [24]. The objective function is as follows,

$$\min_{\mathbf{w}, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{u=1}^U \xi_u, \quad (28)$$

$$\begin{aligned} s.t. \quad & \mathbf{w}^T \delta \Psi_{po}(X^{A,u}, \mathbf{y}^{A,u}, \hat{\mathbf{y}}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V) \\ & \geq \Delta(\mathbf{y}^{A,u}, \hat{\mathbf{y}}^{A,u}) - \xi_u, \\ & \forall \hat{\mathbf{y}}^{A,u} \in \mathcal{Y}^{A,u} \setminus \mathbf{y}^{A,u}, \quad \xi_u \geq 0, \text{ for } u = 1, \dots, U, \end{aligned}$$

where $\delta \Psi_{po}$ is defined as

$$\begin{aligned} \delta \Psi_{po}(X^{A,u}, \mathbf{y}^{A,u}, \hat{\mathbf{y}}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V) \\ = \Psi_{po}(X^{A,u}, \mathbf{y}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V) \\ - \Psi_{po}(X^{A,u}, \hat{\mathbf{y}}^{A,u}; \{X^{B,v}\}_{v=1}^V, \{P^{u,v}\}_{v=1}^V), \quad (29) \end{aligned}$$

\mathbf{w} is the weight vector, C is a parameter to balance between margin and training error, $\mathbf{y}^{A,u}$ is a correct partial order that ranks all correct matches before incorrect matches, and $\hat{\mathbf{y}}^{A,u}$ is an incorrect partial order that violates some of the pairwise relations, e.g. a correct match is ranked after an incorrect match in $\hat{\mathbf{y}}^{A,u}$. The constraints in Eq. (28) force the discriminant score of correct partial order $\mathbf{y}^{A,u}$ to be larger than that of incorrect one $\hat{\mathbf{y}}^{A,u}$ by a margin, which is determined by a loss function $\Delta(\mathbf{y}^{A,u}, \hat{\mathbf{y}}^{A,u})$ and a slack variable ξ_u .

AUC loss function. Many loss functions can be applied in structural SVM. In person re-identification, we choose the ROC Area loss, which is also known as Area Under Curve (AUC) loss. It is computed from the number of swapped pairs,

$$N_{swap} = \{(v, v') : X^{B,v} \succ X^{B,v'} \text{ and} \quad (30)$$

$$\mathbf{w}^T \Phi(X^{A,u}, X^{B,v}; P^{u,v}) < \mathbf{w}^T \Phi(X^{A,u}, X^{B,v'}; P^{u,v'})\},$$

i.e. the number of pairs of samples that are not ranked in a correct order. In the case of partial order ranking, the loss function is

$$\begin{aligned} \Delta(\mathbf{y}^{A,u}, \hat{\mathbf{y}}^{A,u}) &= |N_{swap}| / |S_{X^{A,u}}^+| \cdot |S_{X^{A,u}}^-|, \\ &= \sum_{v, v'} (1 - \hat{y}_{v,v'}^{A,u}) / (2 \cdot |S_{X^{A,u}}^+| \cdot |S_{X^{A,u}}^-|). \end{aligned} \quad (31)$$

We note that there are an exponential number of constraints in Eq. (28) due to the huge dimensionality of $\mathcal{Y}^{A,u}$. Joachims *et al.* [24] showed that the problem could be efficiently solved by a cutting plane algorithm. In our problem, the discriminative model is learned by the structural RankSVM algorithm, and the weight vector \mathbf{w} in our model means how important it is for each term in Eq. (21). In Eq. (21), $\{\alpha_{p_i,k}\}_{k=1,2,3,4}$ correspond to the first four terms based on saliency matching with visual similarity, and $\{\beta_{p_i,k}\}_{k=3,4}$ correspond to the last four terms only depending on saliency matching.

We visualize the learning result of \mathbf{w} in Figure 11, and find that the first four terms in Eq. (21) are heavily weighted in the central part of human body which

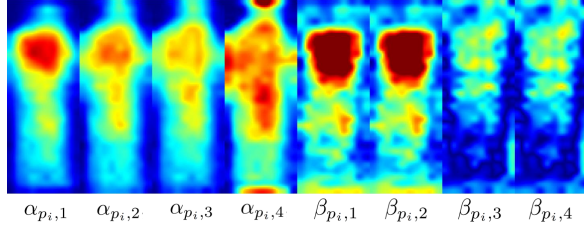


Fig. 11. We Normalize the learnt weight vector \mathbf{w} to a 2-dimensional importance map for different spatial locations. Eight importance maps correspond to $\{\alpha_{p_i,k}\}_{k=1,2,3,4}$ and $\{\beta_{p_i,k}\}_{k=1,2,3,4}$ in Eq. (17).

implies the importance of saliency matching based on visual similarity. $\{\beta_{p_i,k}\}_{k=1,2}$ are not relevant to visual similarity and they correspond to the two cases when $l_{p_i}^A = 1$, i.e. the patches on the query images are salient. It is observed that their weighting maps are highlighted on the upper body, which matches to our observation that salient patches usually appear on the upper body. $\{\beta_{p_i,k}\}_{k=3,4}$ are not relevant to visual similarity either, but they correspond to the cases when $l_{p_i}^A = 0$, i.e. the patches on the query images are not salient. We find that their weights are very low on the whole maps. It means that non-salient patches on query images have little effect on person re-identification if the contribution of visual similarity is not considered.

7.5 Combination with existing approaches

Our approach is complementary to existing approaches. In order to combine existing approaches with the matching score in Eq. (19), the distance between two images can be computed as follows:

$$\text{sim}_{eSalMatch_{opt}}(X^A, X^B) = \sum_i \mu_i \cdot \text{sim}_i(X^A, X^B) - \mu_{Sal} \cdot \text{sim}_{SalMatch_{opt}}(X^A, X^B) \quad (32)$$

where $\mu_i (> 0)$ is the weight for the i th similarity measure, $\mu_{Sal} (> 0)$ the weight for unified saliency matching similarity. sim_i corresponds to the similarity measures using wHSV and MSCR in [13] or LADF [35]. In the experiment, $\{\mu_i\}$ are chosen the same as in [13], [35]. μ_{Sal} is fixed as 1.

8 EXPERIMENTAL RESULTS

We evaluate our approach on two public datasets, i.e. the VIPeR dataset [16], and the CUHK01 dataset [34]. Examples of images in the two datasets are shown in Figure 14. Qualitative results of saliency learning are shown, and quantitative results are reported in standard Cumulated Matching Characteristics (CMC) curves [52].

8.1 Datasets

VIPeR Dataset [16]. The VIPeR dataset ¹ contains images from two cameras, which were placed at many

different locations in an outdoor academic environment. Therefore, the viewpoint changes between cameras are complex. From the time it was publicly available, it has become one of the most challenging person re-identification datasets. It contains 632 pedestrian pairs, each pair contains two images of the same individual seen from different cameras. Most of the image pairs show viewpoint change larger than 90 degree. All images are normalized to 128×48 for experiments.

CUHK01 Dataset [34]. The CUHK01 dataset² was also captured from two camera views in a campus environment. Images in this dataset are of higher resolution and are more suitable to show the effectiveness of saliency matching. It has 971 persons, and each person has two images from camera *A* and the other two from camera *B*. Camera *A* is from a frontal view and camera *B* is from a side view. All images are normalized to 160×60 for evaluations.

The CUHK01 dataset was recently built and contains more images than VIPeR (3884 *vs.* 1264). Both are very challenging datasets for person re-identification because they contain significant variations on viewpoints, poses, and illuminations, and their images are with occlusions and background clutters.

8.2 Evaluation Protocol

Our experiments on both datasets follow the evaluation protocol in [17], i.e. we randomly partition the dataset into two even parts, 50% for training and 50% for testing. Images from camera *A* are used as probe and those from camera *B* as gallery. Each probe image is matched with every image in gallery, and the rank of correct match is obtained. Rank- k matching rate is the expectation of correct match at rank k , and the cumulated values of matching rate at all ranks is recorded as one-trial CMC result. 10 trials of evaluation are conducted to achieve stable statistics, and the expectation is reported. We denote our approach by *SalMatch* for comparison.

8.3 Evaluation on saliency Learning

We investigate the correlation between the human saliency estimated from human perception through user study and that automatically estimated by computation models. The computation models include those design for general image saliency (such as Itti [22] and Hou [21]) and our KNN and One-Class SVM (OCSVM) models specially designed for human saliency. We compute the mean saliency score of each annotated body part, and the Pearson correlation between the automatically estimated saliency and estimation from human perception. Results are shown, The scatter map in Figure 12(a) shows our learned saliency (KNN and OCSVM) has high positive correlations with human perception over the 100 annotated images, while general image saliency (Itti and

1. <http://vision.soe.ucsc.edu/?q=node/178>

2. http://www.ee.cuhk.edu.hk/~xgwang/CUHK_identification.html

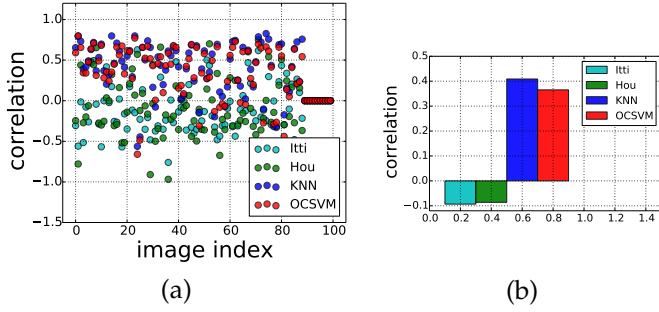


Fig. 12. Correlation between automatically estimated saliency by different approaches (Itti [22], Hou [21], our KNN model and our One-Class SVM (OCSVM) model) and estimation from human perception. (a) Scatter plot of correlations over 100 images. (b) Average correlations.

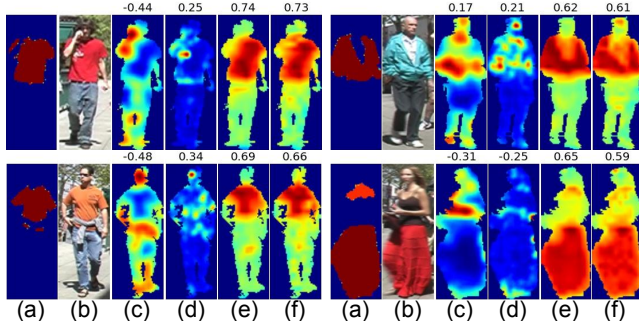


Fig. 13. Examples of estimated saliency map (only body parts are shown). (a) Human saliency estimated from user study. (b) Pedestrian images. (c) and (d) are general image saliency estimated by Itti [22] and Hou [21]. (e) and (f) are human saliency estimated by KNN and OCSVM. Number on top of each saliency map indicates the correlation with human saliency estimated from user study.

Hou) exhibits slight negative correlations. Figure 12(b) shows averaged correlations. Some compared examples are shown in Figure 13. The approaches for general image saliency detection can separate body parts from background. However, the identified body parts may not be effective on recognizing identities.

More interesting results of saliency estimation are shown in Figure 14(a)(b) both on the VIPeR dataset and the CUHK01 dataset. Qualitative results show our saliency learning approach could well approximate human perception and capture important salient regions on human body.

We also quantitatively compare the effectiveness of the saliency estimated from user study and our computation models in person re-identification. We regard the 100 images (of 100 different persons) with saliency estimated from user study as the probe set for evaluation, and images of the corresponding identities in another camera view are included as the gallery set. Bi-directional weighted matching is adopted in testing competing saliency estimation methods, including general image saliency (Itti and Hou), our learned human saliency (SDC_knn and SDC_ocsvm), and saliency estimated from user study (SDC_gt). CMCs are reported in Figure

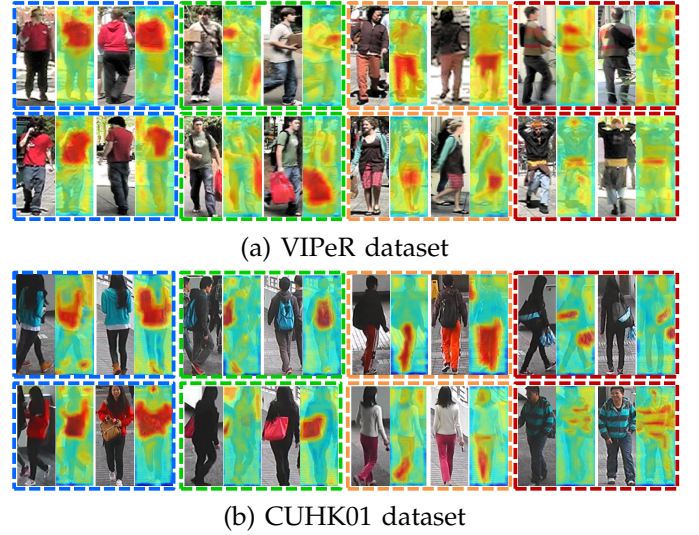


Fig. 14. Examples of saliency matching in our experiments. It shows four types of saliency distributions: saliency in upper body (in blue dashed box), saliency of taking bags (in green dashed box), saliency of lower body (in orange dashed box), and saliency of stripes on human body (in red dashed box). **Best viewed in color.**

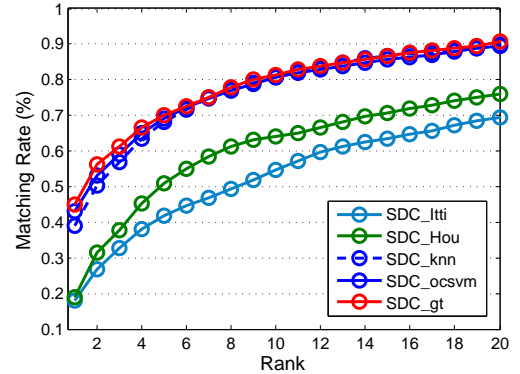


Fig. 15. Bi-directional weighted matching (denoted by *SDC*) using different saliency estimated by different approaches.

15. Results show that the our learned human saliency can well approximate the saliency estimated from user study in person re-identification, while general image saliency significantly degrades the re-identification performance.

8.4 Component-wise Evaluation

The effectiveness of different components in our framework is evaluated. Different settings of component combination are described in Table 1 and their results are shown in Figure 16. DenseFeats in Eq. (2) performs the worst since it directly matches misaligned patches. PatMatch in Eq. (12) performs better by handling misalignment. SDC_knn (SDC_ocsvm) in Eq. (14) improves the performance by incorporating the estimated KNN (One-class SVM) saliency in patch matching. SalMatch_knn (SalMatch_ocsvm) in Eq. (27) formulates person re-identification as saliency

Denotation	Description of component combination in test
<i>DenseFeats</i>	Matching with concatenated patch features
<i>PatMatch</i>	Use patch matching to handle misalignment
<i>SDC_knn</i>	Bi-directional weighted matching (KNN saliency)
<i>SalMatch_knn</i>	Unified saliency matching (KNN saliency)
<i>eSalMatch_knn_1</i>	Combine <i>SalMatch_knn</i> with SDALF [13]
<i>eSalMatch_knn_2</i>	Combine <i>SalMatch_knn</i> with LADF [35]
<i>SDC_ocsvm</i>	Bi-directional weighted matching (OCSVM saliency)
<i>SalMatch_ocsvm</i>	Unified saliency matching (OCSVM saliency)
<i>eSalMatch_ocsvm_1</i>	Combine <i>SalMatch_ocsvm</i> with SDALF [13]
<i>eSalMatch_ocsvm_2</i>	Combine <i>SalMatch_ocsvm</i> with LADF [35]

TABLE 1

Description of all the test settings in components evaluation.
Refer to evaluation results in Figure 16.

matching, and learns matching weights in a supervised way. *eSalMatch_knn_1* (*eSalMatch_ocsvm_1*) in Eq. (32) ensembles SDALF feature matching scores in *SalMatch_knn* (*SalMatch_ocsvm*) matching scores, and *eSalMatch_knn_2* (*eSalMatch_ocsvm_2*) ensembles LADF similarity measures. By combining with either method, the fusion methods outperforms each component, showing that our approach is complementary to other methods. One-class SVM saliency achieves slightly better than its counterpart settings using KNN saliency.

8.5 Comparison with the state-of-the-art

Figure 17 shows significant improvement of *SDC* (unsupervised) compared with existing unsupervised methods, *i.e.* SDALF [13], CPS [11], eBiCov [41], eLDFV [42], and Comb [27] in the VIPeR dataset. For the CUHK01 dataset, we only include the *DenseFeats* and the SDALF in comparison since code or feature representations of the other methods are not available.

Figure 18 compares our supervised saliency matching (*SalMatch* and *eSalMatch*) with ten alternative supervised methods, including seven benchmarking distance metric learning methods, *i.e.* PRDC [58], LMNN-R [12], KISSME [26], LADF [35], PCCA [45], attribute-based PRDC (aPRDC) [37] and LF [46], a boosting approach (ELF) [17], an ensemble of RankSVM (PRSV) [48], and a sparse ranking method (ISR) [36]. Our approach outperforms all these methods. They ignore the domain knowledge on spatial variation caused by misalignment and poses as mentioned in Section 3. Although aPRDC shares a similar spirit as ours in finding unique and inherent appearance, it weights different types global features instead of local patches. Its Rank-1 accuracy is only half of ours. ELF has a low performance since it selects features in the original feature space in which features of different classes are highly correlated. RankSVM is similar to our method in formulating person re-identification as ranking problem. Combined approach *eSalMatch* is not evaluated in CUHK01 dataset because the weights μ_i in Eq. (32) are not carefully tuned for this dataset in SDALF method, and features of this dataset are not available in combining method LADF [35]. Compared with classical metric learning methods (CCA, LMNN, and ITML) based on our *DenseFeats* features in CUHK01 dataset, our approach also has the

best performance, as shown in Figure 18(b). Ours has much better performance because we adopt the discriminative saliency matching strategy for pairwise matching, and the structural SVM incorporates ranking loss in global optimization. This implies the importance of exploiting human saliency matching and its effectiveness in training structural SVM.

9 CONCLUSION AND FUTURE WORK

We propose a novel human saliency learning and matching framework for person re-identification. Adjacency constrained patch matching is applied to build dense correspondence between image pairs to handle misalignment caused by drastic viewpoint change and pose variations. Then K-Nearest Neighbor and One-class SVM approaches are proposed to estimate saliency score for each image patch without using identity labels. User study shows that the automatically estimated human saliency has good correlation with human perception. It is more effective than general image saliency in person re-identification. The estimated saliency can be incorporated into patch matching in both the bi-directional matching scheme and the unified saliency matching framework, and images of the same identity can be recognized by maximizing the saliency matching score. Learning the weights in unified saliency matching framework is formulated as solving a structural RankSVM problem. Experimental results valid the effectiveness of our approach and show superior performances on both the VIPeR and CUHK01 datasets.

The proposed framework can be extended by being integrated with other person re-identification approaches. For example, *DenseFeats* used in this work can be replaced by other more advanced descriptors of characterizing local patches. Patch matching in our framework can be replaced by more sophisticated feature matching techniques [31]. Since saliency information is complementary to appearance, our saliency matching result can be combined with the matching results of existing approaches to boost their performance as shown in Section 7.5.

REFERENCES

- [1] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Susstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Trans. on PAMI*, 34:2274–2282, 2012.
- [2] T. Avraham, I. Gurvich, M. Lindenbaum, and S. Markovitch. Learning implicit transfer for person re-identification. In *Workshops ECCV*, pages 381–390. Springer, 2012.
- [3] T. Avraham and M. Lindenbaum. Learning appearance transfer for person re-identification. In *Person Re-Identification*, pages 231–246. Springer, 2014.
- [4] S. Bak, E. Corvee, F. Brémond, and M. Thonnat. Person re-identification using spatial covariance regions of human body parts. In *Proc. AVSS*, 2010.
- [5] A. Bedagkar-Gala and S. K. Shah. A survey of approaches and trends in person re-identification. *Image and Vision Computing*, 32(4):270–286, 2014.
- [6] A. Borji and L. Itti. Exploiting local and global patch rarities for saliency detection. In *Proc. CVPR*, 2012.

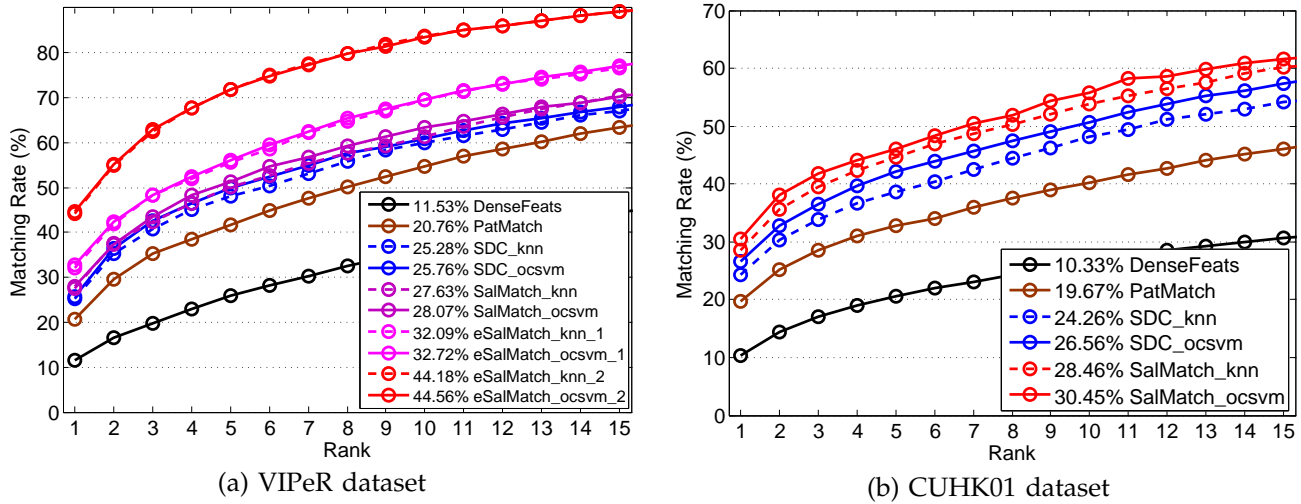


Fig. 16. CMC curves of component-wise evaluation in our approach on the VIPeR and CUHK01 datasets. All the rank-1 accuracies are shown in front of method names.

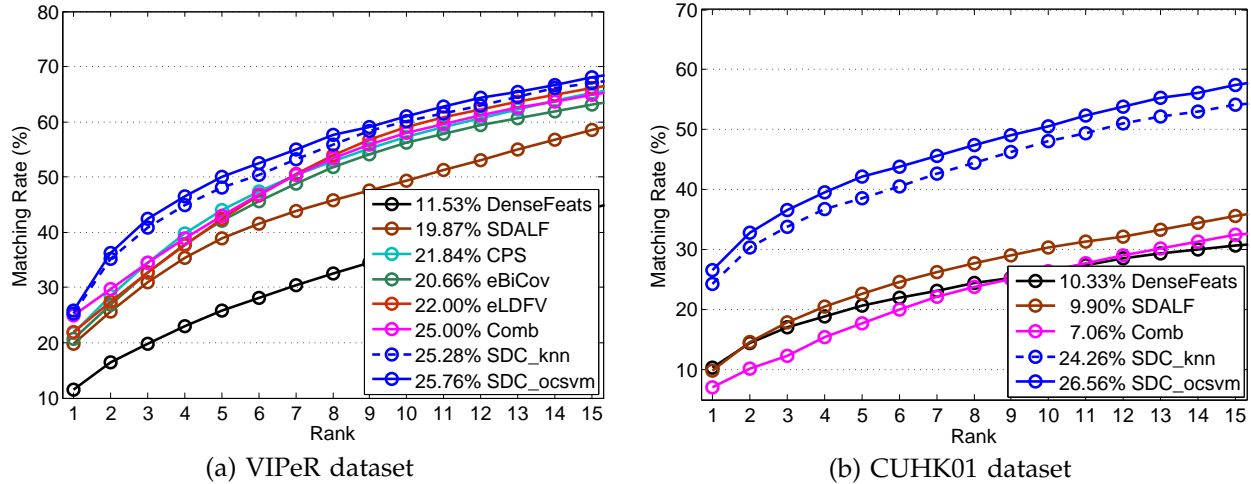


Fig. 17. CMC curves of unsupervised approaches. Rank-1 accuracies are marked in front of method names.

- [7] S. Byers and A. E. Raftery. Nearest-neighbor clutter removal for estimating features in spatial point processes. *Journal of the American Statistical Association*, 93:577–584, 1998.
- [8] B. Carterette and D. Petkova. Learning a ranking from pairwise preferences. In *Proc. ACM SIGIR*, 2006.
- [9] Y. Chen, X. Zhou, and T. Huang. One-class svm for learning in image retrieval. In *Proc. ICIIP*, 2001.
- [10] D. S. Cheng and M. Cristani. Person re-identification by articulated appearance matching. In *Person Re-Identification*, pages 139–160. Springer, 2014.
- [11] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino. Custom pictorial structures for re-identification. In *Proc. BMVC*, 2011.
- [12] M. Dikmen, E. Akbas, T. S. Huang, and N. Ahuja. Pedestrian recognition with a learned metric. In *Proc. ACCV*, 2011.
- [13] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani. Person re-identification by symmetry-driven accumulation of local features. In *Proc. CVPR*, 2010.
- [14] S. Goferman, L. Zelnik-Manor, and A. Tal. Context-aware saliency detection. *IEEE Trans. on PAMI*, 34:1915–1926, 2012.
- [15] S. Gong, M. Cristani, S. Yan, and C. C. Loy. *Person re-identification*. Springer, 2013.
- [16] D. Gray, S. Brennan, and H. Tao. Evaluating appearance models for recognition, reacquisition, and tracking. In *PETS*, 2007.
- [17] D. Gray and H. Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *Proc. ECCV*, 2008.
- [18] K. Heller, K. Svore, A. Keromytis, and S. Stolfo. One class support vector machines for detecting anomalous windows registry accesses. In *Workshop on Data Mining for Computer Security*, 2003.
- [19] M. Hirzer, C. Belezni, M. Kostinger, P. M. Roth, and H. Bischof. Dense appearance modeling and efficient learning of camera transitions for person re-identification. In *Proc. ICIIP*, 2012.
- [20] M. Hirzer, P. M. Roth, M. Kostinger, and H. Bischof. Relaxed pairwise learned metric for person re-identification. In *Proc. ECCV*, 2012.
- [21] X. Hou, J. Harel, and C. Koch. Image signature: Highlighting sparse salient regions. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 34(1):194–201, 2012.
- [22] L. Itti, C. Koch, E. Niebur, et al. A model of saliency-based visual attention for rapid scene analysis. *IEEE Trans. on PAMI*, 20:1254–1259, 1998.
- [23] T. Joachims. A support vector method for multivariate performance measures. In *Proc. ICML*, 2005.
- [24] T. Joachims, T. Finley, and C.-N. J. Yu. Cutting-plane training of structural svms. *Machine Learning*, 77:27–59, 2009.
- [25] T. Judd, K. Ehinger, F. Durand, and A. Torralba. Learning to predict where humans look. In *Proc. ICCV*, 2009.
- [26] M. Kostinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large scale metric learning from equivalence constraints. In *Proc. CVPR*, pages 2288–2295. IEEE, 2012.
- [27] I. Kviatkovsky, A. Adam, and E. Rivlin. Color invariants for person reidentification. *IEEE Trans. on PAMI*, 35(7):1622–1634, 2013.
- [28] R. Layne, T. M. Hospedales, and S. Gong. Towards person identification and re-identification with attributes. In *Workshops ECCV*, pages 402–412. Springer, 2012.
- [29] R. Layne, T. M. Hospedales, S. Gong, et al. Person re-identification by attributes. In *BMVC*, volume 2, page 3, 2012.
- [30] A. Li, L. Liu, and S. Yan. Person re-identification by attribute-assisted clothes appearance. In *Person Re-Identification*, pages 119–138. Springer, 2014.
- [31] H. Li, X. Huang, J. Huang, and S. Zhang. Feature matching

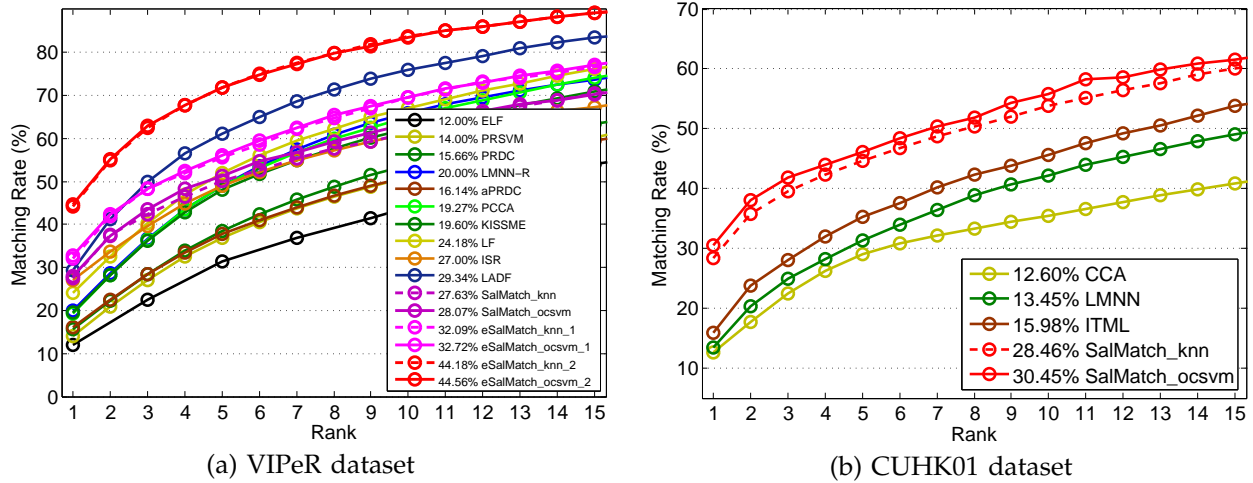


Fig. 18. CMC curves of supervised approaches. Rank-1 accuracies are marked in front of method names.

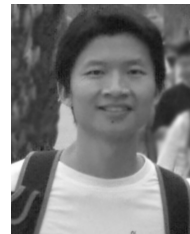
with affine-function transformation models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(12):2407–2422, Dec 2014.

- [32] H. Li and K. N. Ngan. A co-saliency model of image pairs. *IEEE Trans. on Image Processing*, 20:3365–3375, 2011.
- [33] W. Li and X. Wang. Locally aligned feature transforms across views. In *Proc. CVPR*, 2013.
- [34] W. Li, R. Zhao, and X. Wang. Human reidentification with transferred metric learning. In *Proc. ACCV*, 2012.
- [35] Z. Li, S. Chang, F. Liang, T. S. Huang, L. Cao, and J. R. Smith. Learning locally-adaptive decision functions for person verification. In *CVPR*, pages 3610–3617. IEEE, 2013.
- [36] G. Lisanti, I. Masi, A. Bagdanov, and A. Del Bimbo. Person re-identification by iterative re-weighted sparse ranking. *IEEE Trans. on PAMI*, 2014.
- [37] C. Liu, S. Gong, C. C. Loy, and X. Lin. Person re-identification: what features are important? In *Proc. ECCV*, 2012.
- [38] C. Liu, C. C. Loy, S. Gong, and G. Wang. Pop: Person re-identification post-rank optimisation. In *International Conference on Computer Vision*, 2013.
- [39] X. Liu, M. Song, Q. Zhao, D. Tao, C. Chen, and J. Bu. Attribute-restricted latent topic model for person re-identification. *Pattern Recognition*, 45(12):4204–4213, 2012.
- [40] C. C. Loy, C. Liu, and S. Gong. Person re-identification by manifold ranking. In *Proc. ICIP*, volume 20, 2013.
- [41] B. Ma, Y. Su, and F. Jurie. Bicov: a novel image representation for person re-identification and face verification. In *Proc. BMVC*, 2012.
- [42] B. Ma, Y. Su, and F. Jurie. Local descriptors encoded by fisher vectors for person re-identification. 2012.
- [43] B. Ma, Y. Su, and F. Jurie. Discriminative image descriptors for person re-identification. In *Person Re-Identification*, pages 23–42. Springer, 2014.
- [44] B. McFee and G. Lanckriet. Metric learning to rank. In *Proc. ICML*, 2010.
- [45] A. Mignon and F. Jurie. Pcca: A new approach for distance learning from sparse pairwise constraints. In *Proc. CVPR*, 2012.
- [46] S. Pedagadi, J. Orwell, S. Velastin, and B. Boghossian. Local fisher discriminant analysis for pedestrian re-identification. In *Proc. CVPR*, pages 3318–3325. IEEE, 2013.
- [47] B. Prosser, S. Gong, and T. Xiang. Multi-camera matching using bi-directional cumulative brightness transfer functions. In *Proc. BMVC*, volume 8, pages 164–1. Citeseer, 2008.
- [48] B. Prosser, W.-S. Zheng, S. Gong, T. Xiang, and Q. Mary. Person re-identification by support vector ranking. In *Proc. BMVC*, 2010.
- [49] W. Schwartz and L. Davis. Learning discriminative appearance-based models using partial least squares. In *XXII Brazilian Symposium on Computer Graphics and Image Processing*, 2009.
- [50] D. A. Vaquero, R. S. Feris, D. Tran, L. Brown, A. Hampapur, and M. Turk. Attribute-based people search in surveillance environments. In *WACV*, pages 1–8. IEEE, 2009.
- [51] R. Vezzani, D. Baltieri, and R. Cucchiara. People reidentification in surveillance and forensics: A survey. *ACM Computing Surveys (CSUR)*, 46(2):29, 2013.
- [52] X. Wang, G. Doretto, T. Sebastian, J. Rittscher, and P. Tu. Shape and appearance context modeling. In *Proc. ICCV*, 2007.

- [53] Y. Wu, M. Mukunoki, T. Funatomi, M. Minoh, and S. Lao. Optimizing mean reciprocal rank for person re-identification. In *Proc. AVSS*, pages 408–413. IEEE, 2011.
- [54] Y. Xu, L. Lin, W.-S. Zheng, and X. Liu. Human re-identification by matching compositional template with cluster sampling. In *Proc. ICCV*, pages 3152–3159. IEEE, 2013.
- [55] Y. Yue, T. Finley, F. Radlinski, and T. Joachims. A support vector method for optimizing average precision. In *Proc. ACM SIGIR*, 2007.
- [56] R. Zhao, W. Ouyang, and X. Wang. Learning mid-level filters for person re-identification. In *Proc. CVPR*, 2014.
- [57] W.-S. Zheng, S. Gong, and T. Xiang. Associating groups of people. In *Proc. BMVC*, 2009.
- [58] W.-S. Zheng, S. Gong, and T. Xiang. Person re-identification by probabilistic relative distance comparison. In *Proc. CVPR*, 2011.
- [59] W.-S. Zheng, S. Gong, and T. Xiang. Group association: Assisting re-identification by visual context. In *Person Re-Identification*, pages 183–201. Springer, 2014.

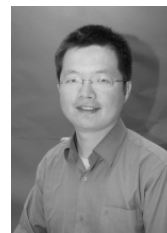


Rui Zhao (S'12) received the B.Eng. degree in Electronic Engineering and Information Science from University of Science and Technology of China in 2010. He is currently a PhD student in the Department of Electronic Engineering at the Chinese University of Hong Kong. His research interests include computer vision, pattern recognition and machine learning.



recognition.

Wanli Ouyang (S'08-M'11) received the B.S. degree in computer science from Xiangtan University, Hunan, China, in 2003. He received the M.S. degree in computer science from the College of Computer Science and Technology, Beijing University of Technology, Beijing, China. He received the PhD degree in the Department of Electronic Engineering, The Chinese University of Hong Kong, where he is now a Research Assistant Professor. His research interests include image processing, computer vision and pattern



Early Career Award in 2012. His research interests include computer vision and machine learning.

Xiaogang Wang (S'03-M'10) received the B.S. degree from University of Science and Technology of China in 2001, the M.S. degree from Chinese University of Hong Kong in 2004, and the PhD degree in Computer Science from Massachusetts Institute of Technology. He is currently an assistant professor in the Department of Electronic Engineering at the Chinese University of Hong Kong. He received the Outstanding Young Researcher Award in Automatic Human Behaviour Analysis in 2011, and the Hong Kong