# A New Patch selection method based on Parsing and Saliency Detection for Person Re-identification

Yixiu Liu<sup>a</sup>, Yunzhou Zhang<sup>a,\*</sup>, Sonya Coleman<sup>b</sup>, Bir Bhanu<sup>c</sup>, and Shuangwei Liu<sup>a</sup>

<sup>a</sup>Northeastern University, College of Information Science and Engineering, Shenyang, China, 110819 <sup>b</sup>Intelligent Systems Research Centre, University of Ulster, Londonderry, UK <sup>c</sup>University of California, Riverside, Riverside, CA-92521

# Abstract

Person re-identification is an important technique towards automatic recognition of a person across non-overlapping cameras. In this paper, a novel patch selection method based on parsing and saliency detection is proposed. The algorithm is divided into two stages. The first stage, primary selection: Deep Decompositional Network (DNN) is adopted to parse a pedestrian image into semantic regions, then sliding window and color matching techniques are proposed to select pedestrian patches and remove background patches. The second stage, secondary selection: saliency detection is utilized to select reliable patches according to saliency map. Finally, PHOG, HSV and SIFT features are extracted from these patches and fused with the global feature LOMO to compensate for the inherent errors of saliency detection. By applying the proposed method on such datasets as VIPeR, PRID2011, CUHK01, CUHK03, PRID 450S and iLIDS-VID, it is found that the proposed descriptor can produce results superior to many state-of-the-art feature representation methods for person identification. Keywords: person re-identification, patch selection, pedestrian parsing, saliency detection, feature fusion

Preprint submitted to Elsevier

<sup>\*</sup>Corresponding author

Email address: zhangyunzhou@mail.neu.edu.cn (Yunzhou Zhang )

# 1 1. Introduction

Person re-identification aims to identify pedestrians in non-overlapping cam-2 eras. It plays a role in a variety of practical applications, such as pedestrian 3 searching, tracking, and analyzing behavior in different camera scenes. Person re-identification makes a significant contribution in reducing time as it can be used to seek a specific person from large amounts of images or videos rather than a human doing this manually. For the above reasons, person re-identification has gained much attention among researchers 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11. However, it remains a challenging problem. A person could undergo significant variations in pose, viewpoint, scale, and illumination when walking through sev-10 eral different cameras. Moreover, background clutter, image blur, and occlusion 11 make the situation even worse. All these problems make intra-person variations 12 even larger than inter-person variations. 13

In this paper, we focus on constructing robust feature representation to solve these problems. Existing methods for feature representation mostly focus on two different aspects: hand-crafted features and deep features 3.

For hand-crafted features, many of them have been developed to achieve 17 precise matching, such as the covariance descriptor based on bio-inspired fea-18 tures (gBiCov) 5, salient color names based color descriptor (SCNCD) 8, and 19 ensemble color model (ECM) [4]. It can be found that these methods have 20 some common problems. They do not remove the background noise, and the 21 attribute types are simplex. Although ECM fuses different color attribute, it has 22 no gradient and other attributes. However, the features we want to construct 23 should have less noise but more diverse attributes. Therefore, a preprocess-24 ing of removing background noise is necessary, and we also consider combining 25 multiple attributes to enhance features. In this paper, Pyramid Histogram of 26 Oriented Gradients (PHOG 12), HSV and Scale Invariant Feature Transform 27 (SIFT 9) features representing gradient, color, and extreme points are fused 28 to complement each other. 29

30

For deep features, they have continuously updated the highest recognition

rate in recent years. A lot of methods are proposed to extract the deep features 31 based on Convolutional Neural Network (CNN). Some of them try to design new 32 CNN frameworks get better deep features, e.g., JointRe-id 13. Some works 33 enhance deep features by fusing with multiple hand-crafted features, e.g., FFN 34 14. Others obtain more discriminating deep features by modifying the loss 35 function in the training process of CNN, e.g., Quadruplet 15. Although each 36 of these method has achieved breakthrough results, we still find their weakness in some practical application scenarios. The problem is that data-driven deep 38 learning cannot play a full role if the samples in the training set are insufficient. 39 So we can see that deep learning methods are usually applied to large-scale 40 person re-identification datasets, such as Market1501 16, DukeMTMC-ReID 41 17, and MSMT17 18. It inspires us to construct a new feature representation 42 to solve the problem of insufficient samples, and the new feature is supposed to 43 improve accuracy more than some deep features. 44

So we consider about picking out valuable patches from images precise fea-15 ture matching. Actually, there are already many researches about local feature 46 exist, much like our idea. Whether you design a local feature (e.g., SDALF 47 **19**) or map an existing local feature space to another space (e.g., LFDA **6**), 48 they all have a same problem: feature drift. The location of most similar patch-49 es in different images changes cross different camera views, and we call this 50 phenomenon feature drift. Some methods try to solve it by saliency features, 51 and get effective improvement, such as SCNCD 8 and SalMatch 20. The fly 52 in the ointment is that they ignore the inherent error of saliency. Therefore, 53 we utilize Graph-Based Visual Saliency (GBVS) 21 to change location-guide 54 feature matching into saliency-guide feature matching, so as to effectively solve 55 the problem of feature drift. In addition, we adopt the strategy combined with 56 global feature Local Maximal Occurrence (LOMO  $(\mathbf{3})$ ) to compensate for the 57 inherent errors caused by saliency detection. 58

In summary, the proposed method makes the following contributions forperson re-identification:

61	(1)	We propose a patch selection method, which can effectively solve the prob-
62		lem of insufficient samples in actual scenarios, and has great significance in
63		engineering applications and has some theoretical value.
64	(2)	In the primary selection, we propose a preprocessing method to remove
65		background noise. We use Deep Decompositional Network (DDN) to divide
66		the picture into semantic regions, and propose sliding window and color
67		matching techniques to remove the background patches.
68	(3)	In the secondary selection, we utilize saliency detection to solve the problem
69		of feature drift and patch unbalance caused by primary selection. It makes
70		us matching features in saliency-guide, rather than location-guide.
71	(4)	We propose a strategy that combines local features with global features to
72		solve the problem of mismatching caused by saliency detection. PHOG,
73		HSV and SIFT features are extracted from the selected patches. LOMO
74		features are extracted from the whole image and fused with them to com-
75		pensate for the inherent error of saliency detection.

The paper is organized as follows. The review of related work is provided in Section 2. The proposed algorithms are described in detail in Section 3. Experimental results using six public benchmark datasets are presented and analyzed in Section 4. Finally, the conclusions are given in Section 5.

# 80 2. Related work

# 81 2.1. Deeply-learned methods for person re-identification

Person re-identification is classified into two categories: single-shot case, and multi-shot case. In general, single-shot person re-identification is required to match a single probe image to a single gallery image. As for multi-shot case, a probe image or images can be matched to frames in the gallery and the matching results can be combined to obtain the result for a video sequence.

In recent years, deep learning has been widely used in image recognition tasks and has made great breakthroughs especially in person re-identification. Yi et al. **13** proposed a method which can simultaneously learn features and

a corresponding similarity metric for person re-identification. Chen et al. [22] 90 presented a novel multi-channel parts-based convolutional neural network (CN-91 N) model that utilized a triplet framework. The CNN model was trained by an 92 improved triplet loss function that assigned the same ID for the closer instances 93 in the learned feature space and assigned a different ID for the farther instances. 94 Furthermore, instead of directly training on the sample images, some methods 95 13, 22, 23, 24 exploited a part or patch-based deep architecture to learn dis-96 criminative feature representations, in local regions of people, with CNNs. For 97 example, Yi et al. 24 split the input image into three rectangular overlapping 98 patches from top to bottom firstly, and then extracted the deep features of each 99 patch through CNN architecture. 100

Through observing the datasets applied by the above methods, we find a rule: the deep learning methods are extremely suitable for large-scale datasets, such as Market1501 [16] and DukeMTMC-ReID [17], and perform well in normal multi-shot datasets, such as CUHK03, but perform relatively poorly in singleshot dataset, such as VIPeR [1]. It inspired us to propose a new method to effectively solve the latter two cases, which is why our approach is only test on datasets such as CUHK03 and VIPeR but not on large-scale datasets.

#### <sup>108</sup> 2.2. Background extraction methods for person re-identification

Background extraction is an important process to improve person re-identification. 109 It separates the target from the background to eliminate the interference of the 110 noisy environment. Based on an improved Random Walks algorithm, Chang et 111 al. 25 proposed an approach that combined the shape prior information and 112 the color seed constraint into the Random Walk formulation, so that each human 113 was divided into several parts where the color features of the HSV histogram 114 and the 1-D RGB signal, along with texture features, were utilized for person 115 re-identification. Le et al., 26 attempted to make a decision on what super-116 pixels belonged to humans and which others belonged to background through 117 the following two techniques: the combination of super-pixels and local saliency 118 information and the combination of super-pixels and pose estimation. 119

Background noise is often ignored in many previous feature representations. In this paper, a new background noise removal strategy is proposed. It is a preprocessing technique of patch primary selection. At first, the pedestrian images are parsed into semantic regions with a Deep Decompositional Network (DDN) [27], such as head, body, arms, and legs. Then pedestrian patches are extracted from the environment using sliding windows and color matching.

#### <sup>126</sup> 2.3. Saliency methods for person re-identification

The saliency of an image carries a lot of potential information that is useful 127 for recognition task. The following methods utilizing saliency are mainly relat-128 ed to human perception in person re-identification. Zhao et al., 28 propose 129 a computational model to estimate the probabilistic saliency map and formu-130 late person re-identification as a saliency matching problem. Saliency matching 131 and patch matching were tightly integrated into a unified structural RankSVM 132 framework. Chen et al., 29 establishes a similarity among patches via fus-133 ing multi-directional saliency after distribution analysis for the consistency of 134 saliency. Le et al., 26 took full advantage of the saliency for keeping super-135 pixels that display a high saliency score (indicating a human) and removing 136 the others (background). In this paper, saliency detection is used for secondary 137 selection, which changes the matching of local features from location-guide to 138 saliency-guide, so as to obtain more reliable patch sequences. 139

# <sup>140</sup> 2.4. Fusion strategy for person re-identification

In this paper, a novel feature representation that combines the global and 141 local features is proposed, which is quite different from other general methods. 142 Most feature-based methods either extract the features from the images directly 143 3.8.30 or use only the local descriptors 1.9.19. Liao 3 designed an efficient 144 feature representation named Local Maximal Occurrence (LOMO), and a sub-145 space and metric learning method known as Cross-view Quadratic Discriminant 146 Analysis (XQDA) 3. Gray and Tao's work 1 proposed an ensemble of invari-147 ant features (EIFs) where the feature representation can effectively handle the 148

variation of human poses/viewpoints and color difference for matching pedestrians observed under different scenes conditions. Every image was divided into
a grid of local patches, and then the color histogram in LAB color space and
SIFT features are extracted for metric learning [28].

The difference from the above works is that we not only fuse many types of features, but also consider the relationship between the global and the local. Considering that our patch selection method has a certain mismatching rate caused by saliency detection, we compensate for it by combining the global features with the local features extracted from the selected patches. The fused feature representation is evaluated with several metric methods which are proven to be effective for person re-identification.

# <sup>160</sup> 3. Technical approach

# <sup>161</sup> 3.1. Structure of the feature representation

The structure of the technical approach consists of three parts: primary selection, secondary selection, and feature fusion. The overall process of the proposed work is shown in Fig. [1].

As can be seen from Fig. 1, parsing and saliency detection are two important 165 techniques for the two patch selection stages, respectively. Throughout the 166 patch selection, our operation unit is patch. Firstly, as a pre-processing, DNN 167 divides the pedestrian's body and background into semantic regions of different 168 colors (3.2.1). It inspires us to propose a patch based sliding window and color 169 matching method to remove the background patches and preserve the pedestrian 170 patches (3.2.2). Afterwards, saliency detection is utilized to get the saliency map 171 (3.3.1), through with we change location-guide feature matching into saliency-172 guide feature matching (3.3.2), and obtain the reliable patch sequences with 173 higher saliency scores (3.3.3). Finally, global and local features are extracted 174 and fused (3.4.1) to obtain complete feature representation, and metric learning 175 is performed to evaluate it (3.4.2). 176



Figure 1: The architecture of the feature representation. It consists mainly of three parts, *i.e.*, primary selection, secondary selection, and feature fusion. (1). For each pair of pictures, we split them into patches and parse them into semantic regions with DDN which will be described in detail in 3.2. The overlap rate is computed to remove the background patches. The threshold of overlap rate is set to 25%. (2). The patches are further selected by saliency detection. The patch sequences with higher saliency scores are obtained. e.g., the most reliable patches that have higher saliency scores are pained in red. (3). We extract the PHOG, HSV histogram and SIFT features from the selected patches and fuse them with global LOMO features.

### 177 3.2. Primary selection

#### 178 3.2.1. Semi-supervised DDN

It is not feasible to fine-tune DDN model directly on person re-identification datasets, because there are no ground truth of label maps. In other words, person re-identification datasets have no label for DDN model. So we modify it into a semi-supervised DDN model, and the training loss function is defined as

$$L = \sum_{x^l} C(y^l, y^l) + \lambda \sum_{x^u} E(y^u), \tag{1}$$

The first term in Eq. [] is the loss function trained on the labelled parsing dataset. The second term is the loss function trained on the unlabelled person re-identification dataset. Before that, let's review the original DDN.

Fig. 2 illustrates the architecture of the DDN which directly maps low-level
visual features to the label map of body parts. The input is the feature vector,
while the output is a set of label maps of body parts. This architecture is
utilized for pedestrian parsing, and mainly consists of one down-sampling layer,
two occlusion estimation layers, two completion layers, and two decomposition layers.



Figure 2: DDN architecture, which achieves parsing and subtraction in a unified deep network.

191

The input x is down-sampled to  $x^d$ . Otherwise, x is mapped into a binary occlusion mask  $x^o \in [0, 1]^n$  through the weight matrices  $w^{o_1}$ ,  $w^{o_2}$  and the biases  $b^{o_1}$ ,  $b^{o_2}$ . To reduce the number of parameters in the network,  $x^o$  and  $x^d$  are set to the same size. If the *i*-th element of the feature is occluded,  $x_i^o$  is set to 0, otherwise  $x_i^o = 1$ . The binary occlusion mask  $x^o$  is denoted as

$$x^{0} = \tau(W^{o_{2}}\rho(W^{o_{1}}x + b^{o_{1}}) + b^{o_{2}}), \qquad (2)$$

where the function  $\tau(x) = 1/(1 + \exp(-x))$ . For the first layer of occlusion estimation, the rectified linear function **31**  $\rho(x) = \max(0, x)$  is utilized as the activation function and we use a sigmoid function as the activation function inthe second layer.

In the architecture of the DDN, the input of the completion layers which are 201 modeled as the denoising autoencoder (DAE) 32 is the element-wise product 202 of  $x^{o}$  and  $x^{d}$ . While the output is the completed feature vector  $x^{c}$  via the 203 weight matrices  $W^{c_1}$ ,  $W^{c_2}$ ,  $W^{c_{1'}}$ ,  $W^{c_{2'}}$ , and the biases  $b^{c_1}$ ,  $b^{c_2}$ ,  $u^{c_1}$ ,  $u^{c_2}$ . W'204 is the transpose of W. Through projecting high dimensional data into a low 205 dimensional space, the encoders  $W^{c_1}$  and  $W^{c_2}$  find the compact representation 206 of noisy data. The encoders  $W^{c_{1'}}$  and  $W^{c_{2'}}$  are used to reconstruct the data. 207 We reconstruct  $x^c$  with  $x^o$  and  $x^d$ . The reconstruction process is as follows. 208

$$z = \rho(W^{c_2}\rho(W^{c_1}(x^o \odot x^d) + b^{c_1}) + b^{c_2}), \tag{3}$$

where  $\odot$  represents the element-wise product, and z denotes the compact representation. According to Eq. 3, we can get

$$x^{c} = \rho(W^{c_{1'}}\rho(W^{c_{2'}}z + u^{c_{2}}) + u^{c_{1}}), \tag{4}$$

At the back end of DDN, the completed feature  $x^c$  is decomposed into several label maps from  $y_1$  to  $y_M$  through the corresponding weight matrices  $W^{t_1}$ ,  $W_1^{t_2},..., W_M^{t_2}$ , and biases  $b^{t_1}, b_1^{t_2},..., b_M^{t_2}$ . We denote the label map  $y_i \in [0,1]^n$  as

$$y_i = \tau(W_i^{t_2}\rho(W^{t_1}x^c + b^{t_1}) + b_i^{t_2})$$
(5)

215 So the loss function for labelled parsing dataset becomes

$$\sum_{x^{l}} C(y^{l}, \overset{\wedge}{y^{l}}) = || \overset{\wedge}{Y^{l}} - Y^{l} ||_{F}^{2}$$
(6)

<sup>216</sup> Where  $Y^l = \{y_i^l\}$  and  $\overset{\wedge}{Y^l} = \{\overset{\wedge}{y^l}\}$  are the set of outputs and the set of ground <sup>217</sup> truth labels.

Now we use the current DDN to train the unlabelled person re-identification dataset. The training follows the hypothesis of low-density separation 33. Specifically, the object of our training is to make the probability that the output



Figure 3: The test results of the VIPeR dataset for person re-identification with DDN.

tends to a class close to 1, and the sum of the probabilities toward other classes
tend to be zero. We define the loss as an entropy

$$\sum_{x^{u}} E(y^{u}) = -\sum_{i=1}^{N} y_{i}^{u} \ln(y_{i}^{u}), \tag{7}$$

Where N denotes the number of samples and  $y_i^u$  is the output. Finally we get

<sup>224</sup> the semi-supervised DDN loss function

$$L = || \stackrel{\wedge}{Y^l} - Y^l ||_F^2 - \sum_{i=1}^N y_i^u \ln(y_i^u) \tag{8}$$

225 3.2.2. Background noise removal

After segmenting the images of pedestrians into a set of semantic regions, we propose a method based on the use of sliding windows and color matching to remove the cluttered environment around the pedestrians. At first, every image is divided into a grid of local patches, and then the background is masked through computing the overlap rate between the mask and patches. This mask is preset, such as the pedestrian's upper body is green and the background is dark blue. The whole process is shown in Fig. 4

233

Every image is divided into patches of size  $10 \times 10$ , with a step size of 5 pixels.



Figure 4: The process of masking the background based on sliding windows and color match.

<sup>234</sup> To determine if a patch will be masked, we apply with following equation:

$$c(P_{ij}) = \frac{u(M - P_{ij})}{x_p * y_p},$$
(9)

where  $P_{ij}$  indicates the patch at the *i*-th row and *j*-th col of the image,  $i, j \in N_+$ , 235  $\{i, j | i \leq m, j \leq n\}$ .  $c(P_{ij})$  denotes the overlapping rate between the sliding 236 mask M and the  $P_{ij}$  and u(x) is indicated as the number of non-zero elements in 23 matrix x.  $x_p$  and  $y_p$  represent the number of patches in the horizontal and ver-238 tical direction, respectively. The patches for which  $c(P_{ij}) \le 25\%$  are reserved, 239 whereas others are masked. Because the background of the same pedestrian 240 often changes under different cameras, background noise removal focuses fea-241 tures on pedestrian patches by removing background patches, making feature 242 matching more accurate. We define all the reserved patches of each image as 243 the set S1. 244

# 245 3.3. Secondary selection

The primary selection may cause two problems, one is feature drift and the other is patch unbalance. The location of most similar patches in different images changes cross different camera views. The number of pedestrian patches selected from different images may not be same, which may lead to different feature lengths. Person saliency is distinctive and reliable in pedestrian matching across disjoint camera views. If the patches of two images from the same person are matched, the saliency values of these patches should be similar to each

<sup>253</sup> other, regardless of their location. In addition, the number of patches is easily

- <sup>254</sup> controlled by saliency scores, so as to keep the features consistent in length.
- 255 3.3.1. Saliency detection
- Based on human focus of attention [28], salient regions are defined with
  the following properties: 1) making the pedestrian more distinctive than other distractors; 2) being reliable to search for the same pedestrian across different



Figure 5: Silent region could be the part of the human body or the decorations the person carries. The salient regions are circled with the yellow dotted lines.

258

camera views. Compared with the abstract features, it's easier for a human to identify the same person, because if the salient region occurs in one camera view, it usually remains salient in another camera view. For example, in Fig. a human would easily identify that there is a red bag on the shoulder of person p1, p2 carries a yellow bag, p3 has a red umbrella in his hand while p4 holds a green parcel in his hand.

A reliable approach to map the salient regions is saliency learning 28. It divides pedestrians into different parts and manually merges super-pixels that are coherent in appearance. Then the segmented body part is randomly selected and presented to a labeler. The labeler is allowed to select the most likely image from the list based on visual perception. However, this method requires a significant amount of man hours, so it is impractical for large datasets. In this paper, GBVS 21 is employed to automatically detect the salient regions.  $_{272}$  Moreover, to reduce the huge cost of matching time, we select only 25 patches

 $_{\rm 273}\,$  whose saliency scores are relatively higher than others. This number is the

empirical result of the compromise between computation time and matching accuracy.



Figure 6: Illustration of saliency detection with the GBVS algorithm and the saliency map of the pedestrian image is shown. Best viewed in color.

275

As we can see from Fig. 6 that the salient region is detected by the GBVS 276 algorithm that computes bottom-up saliency maps which show a remarkable 277 consistency with the attentional deployment of human subjects. In many cases, 278 different persons from different camera views have different spatial distribution, 279 whereas the salient region of the same pedestrian under different camera views is 280 discriminative from others. For example, the salient region in (a1) is a backpack. 281 The similar salient region also exists in  $(a_2)$ , so  $(a_2)$  is the correct match of  $(a_1)$ . 282 There is a green bag hanging on the pedestrian's arm in (a3). The yellow bag 283 on the shoulder of the woman in (a4) is very eye-catching. While the woman 284 in (a5) holds a white paper in his hand. They are all the incorrect matches of 285 (a1). For the same reason, (b2) is the correct match of (b1). (b3), (b4), (b5) are 286 the incorrect matches of (b1). 287

#### 288 3.3.2. Saliency-guide matching

We hope to match the features of similar patches in different images, but in fact, due to the change of pose and views under different cameras, they are offset in position, and even some patch features may shift from pedestrian to background. Now we change the feature matching from location-guide to saliency-guide, which effectively solves the problem of mismatching caused by feature drift.

At first, the image is constructed as a Gaussian pyramid to extract multiscale features in the down-sampling process.

$$R(\sigma) = I(x, y) \otimes G(x, y, \sigma), \tag{10}$$

297

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{\left(-\frac{x^2 + y^2}{2\sigma^2}\right)},$$
(11)

where  $R(\sigma)$  is the initial feature map using the GBVS model, I(x, y) represents the image,  $G(x, y, \sigma)$  denotes Gaussian pyramid,  $\sigma$  is the scale factor or bandwidth of Gaussian pyramid and  $\otimes$  in Eq. 10 denotes the convolution operator. Secondly, the activation maps are formed using the feature maps, and the most important thing is to construct the Markov matrix. We assume that the scale of the feature graph is constant. In other words, we ignore the scale  $\sigma$ . We then define the dissimilarity of R(x, y) and R(p, q) as

$$d((i,j)||(p,q)) \triangleq \left| \log \frac{R(i,j)}{R(p,q)} \right|,$$
(12)

where R(x, y) and R(p, q) represent the feature value of the pixels at (i, j) and (p, q), respectively. We obtain the fully-connected directed graph  $G_A$  through connecting every node of the lattice R, labelled with the indices (i, j) or (p, q). The directed edge from node (i, j) to node (p, q) will be assigned a weight

$$w_1((i,j),(p,q)) \stackrel{\Delta}{=} d((i,j)||(p,q)) \cdot F(i-p,j-q),$$
(13)

309

$$F(a,b) \stackrel{\Delta}{=} \exp\left(-\frac{a^2 + b^2}{2\sigma^2}\right),\tag{14}$$

where  $\sigma$  is a constant which denotes the free parameter. The Markov chain is defined on directed graph  $G_A$ . We normalize the weights on the edges of  $G_A$ to be 1. Now the stationarity of the Markov chain is utilized to obtain the probability that the state node transforms to another, thereby estimating the saliency of the directed graph and obtaining the saliency map A.

Finally, we normalize the saliency map A, and construct the directed graph  $G_N$ . We redefine a Markov chain on  $G_N$ , and introduce an edge from (i, j) to 317 (p,q) with weight:

$$w_2((i,j),(p,q)) \stackrel{\Delta}{=} A(p,q) \cdot F(i-p,j-q), \tag{15}$$

where A denotes the final saliency map; every element inside represents the saliency value of the pixel in this position. The size of A is the same as the original image. Every image is divided into patches of size  $10 \times 10$  with a step size of 5 pixels, and the patches which have the higher saliency value are selected by

$$s\left(\left\{p_A\right\}(i,j)\right) = average\left(\left\{p_A\right\}(i,j)\right),\tag{16}$$

where  $p_A(i, j)$  denotes the patch at the *i*-th row and *j*-th column of A,  $s(p_A(i, j))$ is the average saliency value of  $p_A(i, j)$ . We use 0.6 as the empirical value of  $s(p_A(i, j))$ . The patches are reserved corresponding to the original image where  $s(p_A(i, j)) \ge 0.6$ , while others are removed. We define all the reserved patches of each image as the set S2.

#### 328 3.3.3. Aligned patch sequences

After primary selection and secondary selection, we obtain corresponding 329 patch sets S1 and S2 respectively. Now we define their intersection  $S = S1 \cap S2$ 330 as a set of reliable patches. Due to the different views under different cameras, 331 the proportions of pedestrian and background are also different. Some images 332 have more pedestrian patches than background patches, while others are the 333 opposite. This problem of patch unbalance results in a different number of 334 reliable patches selected per image, and correspondingly different lengths of 335 extracted features. 336

In order to ensure that the dimension of the local features extracted from each image is the same, 25 patches are selected from each image from camera A which have a relatively high saliency value within the set S. Using the priori saliency spatial distribution of these patches, we find 25 patches corresponding to the previous 25 patches from each image from camera B with the nearest neighbor classifier for saliency.

Now the similarity of saliency between the patch pairs for different images across disjoint camera views is defined as

$$sim_{saliency}(P^{A,u}, P^{B,v}) = \exp(-\frac{d(p_i^{A,u}, p_j^{B,v})^2}{2\sigma_d^2})$$
 (17)

Saliency patches of a pedestrian image are represented as  $P^{A,u} = \{p_i^{A,u} | i = 1, 2..., 25\}$ , where (A, u) denotes the *u*-th image under camera A, *i* denotes the position of the patch in this image, and  $p_i^{A,u}$  is the saliency vector of the patch.  $d(\cdot)$  is the Euclidean distance, and  $\sigma_d$  is a bandwidth parameter. Finally, we get the corresponding patches of images from camera B.

$$I^{B,u} = find(\min(sim_{saliency}(P^{A,u}, P^{B,v})))$$
(18)

350 The special form is

$$I_{i}^{B,u} = find(\min(\exp(-\frac{d(p_{i}^{A,u}, p_{j}^{B,v})^{2}}{2\sigma_{d}^{2}}))),$$
(19)

where  $find(\cdot)$  denotes finding the indexes of patches of an image from camera B according to the saliency matching with the patches of the image from camera A.  $I_i^{B,u}$  is an element of  $I^{B,u}$  which denotes the indexes set as mentioned above,  $i \in \{1, 2...25\}$ .

#### 355 3.4. Feature fusion and metric learning

In order to overcome the shortcomings of either of the methods and take advantage of them, the global features and local descriptors are fused in the process of metric learning, so that we can clearly separate the different pedestrians.

#### 359 3.4.1. Feature extraction and fusion

The features we fuse consist of one global feature (LOMO) and three lo-360 cal features (PHOG, HSV, SIFT). We adopt the strategy of combining global 361 feature and local features to compensate for inherent errors of saliency detec-362 tion which may result in mismatching. Specifically, PHOG contains oriented 363 gradient, HSV reflects color distribution, and SIFT captures extreme points in 364 images. As about the global feature (LOMO), although it also contains some 365 color information, it reflects the color distribution of the whole image, using im-366 age pair matching rather than saliency-guide patch pair matching. In a word, 367 they complement each other without redundancy. 368

The LOMO algorithm analyzes the horizontal occurrence of local features, 369 and maximizes the occurrence to make a stable representation against view-370 point changes. Besides, to handle illumination variations, the Retinex trans-371 form **3** and a scale invariant texture operator are applied. To make person 372 re-identification easier than using original images, we apply the HSV color his-373 togram to extract features which has  $8 \times 8 \times 8$  bins = 512 dimensions. The Scale 374 Invariant Local Ternary Pattern (SILTP) 34 descriptor is also extracted for re-375 ducing the impact of illumination invariant. SILTP is an improved operator over 376 the well-known Local Binary Pattern (LBP) 35. We utilize sliding windows 371 with a size of  $10 \times 10$  pixels and an overlapping step of 5 pixels to locate local 378 patches in  $128 \times 48$  pixel images. Two scales of SILTP histograms (SILTP<sub>4,3</sub> 379 and  $SILTP_{4.5}^{0.3}$ ) are extracted, and the dimension of SILTP is  $3^4 \times 2 = 81$ . A 380 three-scale pyramid representation is built for utilizing the multi-scale informa-381 tion, which down-samples the original  $128 \times 48$  image by two  $2 \times 2$  local average 382 pooling operations and then repeats the above feature extraction procedure. So 383 the final feature has  $(8 \times 8 \times 8 \text{ color bins} + 3^4 \times 2 \text{ SILTP bins}) \times (24 + 11 + 5)$ 384 horizontal groups) = 26960 dimensions. 385

On the other hand, the PHOG, HSV histogram and SIFT features are extracted from every selected patch. PHOG is the Pyramid Histogram of Oriented Gradient, which is an effective descriptor for classification; it is the splicing of

the HOG features at different scales. In this work, the number of layers of 389 pyramids is L = 3, and the number of bins of gradient division is n = 8. The 390 dimension of PHOG features is  $(1 + 4 + 16 + 64) \times 8 = 680$ . Color histogram 391 is a significant descriptor that performs outstandingly in the recognition task. 392 To obtain the HSV histogram features, the RGB image is converted to a HSV 393 image at first. The dimension of the HSV histogram feature is  $8 \times 8 \times 8 = 512$ . 394 Besides, we extract the SIFT features that has 128 dimensions. The schematic 395 representation of feature extraction is shown in Fig. 1 396

After finishing the feature extraction, we obtain features with the dimen-397 sion of  $26960 + (680 + 512 + 128) \times 20 = 53360$ . Before concatenating them, we 398 fuse them based on metric learning. Due to the diversity of our features and 399 the complexity of the processing process, we count the time consumption of 400 feature extraction and algorithm execution. First, we count the time of fea-401 ture extraction and fusion on the VIPeR dataset, and the average time for 402 each image was 46.6ms. The experiment is repeated 10 times and averaged (i7-403 6700 CPU, 2.60GHz, Matlab, Windows). Then we perform our algorithm on 404 CUHK03 dataset including semi-supervised DDN training, which takes a total 405 2586,463 ms  $\approx 43$  m 6s, about half the time of FFN 14 (Titan xp, 12GB video 406 memory, GPU, Linux). 407

#### 408 3.4.2. Metric learning

We define  $dist_{i,j}$  as the distance between the features  $x_i$  and  $x_j$  cross different camera views.

$$dist^{2}(x_{i}, x_{j}) = ||x_{i} - x_{j}||_{2}^{2}$$
  
=  $w_{1} \cdot dist^{2}_{ij,1} + w_{2} \cdot dist^{2}_{ij,2} + ... + w_{d} \cdot dist^{2}_{ij,d}$  (20)  
= $(x_{i} - x_{j})^{\mathrm{T}}W(x_{i} - x_{j}),$ 

where  $w_i \ge 0$ , W = diag(w) is a diagonal matrix, and  $(W)_{ii} = w_i$ . W can be determined by learning. d denotes the dimension of the feature which is equal to 53360 in this paper. We replace W with a common semi-definite symmetric matrix M, so we get Mahalanobis distance.

$$dist_{mah}(x_i, x_j) = (x_i, x_j)^{\mathrm{T}} M(x_i, x_j) = ||x_i - x_j||_M^2,$$
(21)

M denotes the metric matrix which is obtained through metric learning. Note that M is the semi-definite symmetric matrix. M is directly embedded into the evaluation of the neighbor classifier, and we obtain M through optimizing the performance of the evaluation. Now we discuss the acquisition of M with the Neighbourhood Component Analysis (NCA) as an example.

Neighbour classifiers use the majority voting method when making a decision. Each sample in the neighbourhood casts one vote, and the samples outside the field casts zero votes. For sample  $x_j$ , the probability of its effect on  $x_i$  classification is

$$p_{i,j} = \frac{\exp(-||x_i - x_j||_M^2)}{\sum_l \exp(-||x_i - x_l||_M^2)},$$
(22)

where l is the number of the samples. As can be seen from Eq. 22,  $p_{i,j}$  is the largest when i = j. If we recognize the maximum accuracy as an optimal object, the accuracy based on leave-one-out (LOO) is computed as follows

$$p_i = \sum_{j \in \Omega_i} p_{ij},\tag{23}$$

<sup>427</sup> where  $\Omega_i$  represents the set of subscripts that belong to the same class as  $x_i$ . <sup>428</sup> The accuracy for the entire sample set is

$$\sum_{i=1}^{m} p_i = \sum_{i=1}^{m} \sum_{j \in \Omega_i} p_{ij} \tag{24}$$

Then we substitute the Eq. 22 into the 24 and make  $M = PP^{T}$ , we get the NCA optimal object

$$\min_{P} = 1 - \sum_{i=1}^{m} \sum_{j \in \Omega_{i}} \frac{\exp(-||P^{\mathrm{T}}x_{i} - P^{\mathrm{T}}x_{j}||_{2}^{2})}{\sum_{l} \exp(-||P^{\mathrm{T}}x_{i} - P^{\mathrm{T}}x_{l}||_{2}^{2})}$$
(25)

Through solving Eq. 25, we obtain the metric matrix M that maximizes the accuracy of the neighbour classifier.

Finally, we get Cumulative Match Characteristic (CMC) curves of person re-identification. Using several different metric methods, we do experiment on different datasets to prove that our proposed method is more effective than many state-of-the-art methods.

# 437 4. Experimental results

There are several existing challenging benchmark datasets for person reidentification. In this work, we perform experiments using six datasets, VIPeR [1], PRID2011 [37], CUHK01 [38], CUHK03 [36], PRID 450S [39], iLIDS-VID [40], which are public benchmarks available to conduct experiments. We emphasize that our approach can effectively solve the problem of insufficient samples in actual scenarios. This is why our method has not been tested in large datasets such as Market1501 [16], DukeMTMC-ReID [17], and MSMT17 [18].

445 4.1. Parameters and implementation details

The parameter settings in this paper are shown in the Table 1. In addition, we perform fine tuning on the basis of original DDN, so the initialization
of parameters and bias are the result of previous training. Two scales of center/surround Retinex is used for image preprocessing when LOMO features are
extracted. For all the experiments, we repeat the procedure 10 times to calculate an average performance.

#### Table 1:

#### Parameter settings

Paranet	ers	Values	Descriptions						
$s_p$		$10 \times 10$	the size of patches						
$s_w$		$10 \times 10$	10 sliding windows						
$c_{thr}$		0.25	0.25 the threshold of overlapping rate $c(P_{i,j})$						
$n_P$		25	the number of patches we selected based on saliency detecti						
$s_{thr}$	$s_{thr}$ 0.6 the threshold of saliency value								
σ		0.5	the scale factor of Eq. 11						
$\sigma_d$	$\sigma_d$ 0.5 bandwidth parameter of Eq. 17								

451

#### 452 4.2. Comparison with state-of-the-art methods

We perform a number of experiments and the results show that the proposed algorithm achieves better performance than many of the existing methods. In



Figure 7: CMC curves of VIPeR, PRID2011, CUHK01, CUHK03, PRID 450S, iLIDS-VID datasets.

order to demonstrate the advantages of our method in the case of insufficient
samples, we also compared the with many deep learning methods, which are
listed separately in the tables. Fig. 7 shows the CMC curves for different
methods on every dataset. The red solid lines represent the results of our
algorithm. It can be seen from Fig. 7 that our method has the highest matching
rate.

# Table 2:

Method	rank=1	rank=10	rank=20	Reference
Ours	56.83	92.03	97.27	Proposed
FFN <b>14</b>	51.1	91.4	96.9	2016 WACV
EBb [41]	51.9	84.8	90.2	2018 CVPR
MLCS 42	34.58	80.59	90.43	2017 TCSVT
LDCA [11]	38.08	73.52	82.91	2017  CVPR
SCSP 43	53.5	90.2	96.6	2016  CVPR
LSSL 44	47.8	87.6	94.2	2016 AAAI
LOMO+XQDA 3	40.00	80.51	91.08	2015  CVPR
SCNCD 8	37.80	81.20	90.40	$2014 \ \mathrm{ECCV}$
kBiCov 5	31.11	70.71	82.45	2014 IVC
SalMatch 20	30.16	65.54	79.15	2013 ICCV
Mid-level Filter <b>9</b>	39.11	65.95	79.87	2014  CVPR
SSCDL 45	25.60	68.10	83.60	2014  CVPR
MtMCML 46	28.83	75.82	88.51	2014 TIP
ColorInv 35	24.21	57.09	69.65	2013 TPAMI
LF <b>6</b>	24.18	67.12	82.00	2013  CVPR

TOP r RANK MATCHING ACCURACY (%) ON VIPeR DATASET.

# 461 4.2.1. Experiments on VIPeR

The VIPeR dataset contains two cameras, each of which captures one image 462 per person. It also provides the viewpoint angle for each image. It has been 463 used by many researchers and is still one of the most challenging datasets. 464 The VIPeR dataset contains 632 pedestrian image pairs taken from arbitrary 465 viewpoints under varying illumination conditions. It is randomly split into two 466 subsets containing the same number of pictures for training and test respectively. 467 We evaluated the proposed algorithm and several state-of-the-art algorithm-468 s, Fig. 7(a) shows the results of the comparisons through CMC curves on the 469 VIPeR dataset. The cumulative matching scores (%) at rank 1, 10, and 20 are 470 listed in Table 2. From Table 2 it can be seen that our method is superior to 471 all compared state-of-the-arts, surpassing the  $2^{nd}$  best method by 3.33% (56.83) 472 53.5) in Rank-1, 0.63% (92.03 - 91.4) in Rank-10, and 0.37% (97.27 - 96.9) in 473 Rank-20. Compared to eliminating background-bias (EBb) method, our method 474 improves the rank-1 by 4.93%, rank-10 by 7.23%, and rank-20 by 7.07%. It in-475 dicates the superiority of primary patch selection by background noise removal. 476 Compared to the deep learning method FFN, our method improves the Rank-1 477 by 5.73%. This indicates that in the single-shot case, our feature fusion strategy 478 is more effectively than FFN that fuses deeply learning features with multiple 479 hand-crafted features. 480

# 481 4.2.2. Experiments on PRID2011

The PRID2011 dataset consists of images extracted from multiple person trajectories recorded from two different, static surveillance cameras. Images from these cameras contain a viewpoint change and a stark difference in illumination, background and camera characteristics. The PRID dataset has 385 trajectories from camera A and 749 trajectories from camera B. Among them, only 200 people appear in both cameras.

Fig. 7(b) shows the results of the comparisons through CMC curves on the PRID2011 dataset. The cumulative matching scores (%) at rank 1, 10, and 20 are listed in Table 3. From Table 3 it can be seen that our method is

#### Table 3:

Method	rank=1	rank=10	rank=20	Reference
Ours	78.3	92.6	97.5	Proposed
VGG-GRU+TP_Mean $[47]$	75.1	97.5	99.5	2017 ICIC
LBP&Color+RFA-Net+RankSVM [48]	58.2	93.4	97.9	$2017 \ \text{ECCV}$
LBP&Color+RFA-Net+Cosine [48]	54.9	93.7	98.4	$2017 \ \text{ECCV}$
RC [49]	70.9	82.7	87.3	2018  CVPR
DVDL [50]	40.6	77.8	85.6	$2015 \ ICCV$
STFV3D+KISSME [2]	64.1	89.9	92.0	2012  CVPR
TDL <b>51</b>	56.7	87.6	93.4	2016  CVPR

## TOP r RANK MATCHING ACCURACY (%) ON PRID2011 DATASET.

<sup>491</sup> superior to all compared state-of-the-arts in Rank-1. It surpasses the 2<sup>nd</sup> best
<sup>492</sup> VGG-GRU+TP\_Mean by 3.2% (78.3 - 75.1) in Rank-1. Although it is 4.9%
<sup>493</sup> (97.5-92.6) and 2.0% (99.5-97.5) lower than VGG-GRU+TP\_Mean in rank-10
<sup>494</sup> and rank-20 respectively, it is not inferior to the suboptimal LBP&Color+RFA<sup>495</sup> Net+Cosine. As can be seen that compared with multiple hits, our method has
<sup>496</sup> a obvious advantage in accuracy of one hit.

#### 497 4.2.3. Experiments on CUHK01

The CUHK01 dataset contains two images for every identity from each camera. This dataset has one pair of disjoint cameras and the image quality of this dataset is relatively good. It contains 971 persons captured from two camera views. Camera A captures the frontal or back views of pedestrians while camera B captures them in a side view.

The CMC curves of comparison with other algorithms are described in Fig. 7(c) All the corresponding data are recorded in Table 4. The proposed method achieved 83.2% at rank-1, which slightly outperforms the second one EBb with an improvement of 0.7% (83.2 - 82.5). The proposed method is comparable to the most advanced algorithms EBb at rank-15 and rank-20, which achieved accuracy at 98.4% and 98.8%, respectively. Compared to SalMatch, we far surpassed it in all the results. It shows that using background noise removal,

#### Table 4:

Method	rank=1	rank=10	rank=15	rank=20	Reference
Ours	83.2	97.1	98.4	98.8	Proposed
Quadruplet 15	62.6	86.0	88.9	89.8	2017  CVPR
MCPB_CNN 22	53.7	91.0	95.4	96.3	2016  CVPR
JointRe-id 13	47.5	80.0	86.8	87.9	2015  CVPR
EBb [41]	82.5	98.2	98.7	99.0	2018 CVPR
LSSCDK 52	66.0	90.0	93.3	95.0	2016  CVPR
Kernel X-CRC 53	61.2	87.3	91.2	93.2	2019 JVCIR
CVPDL 54	59.5	89.7	91.7	93.1	2015 ICOAI
Ensembles 55	51.9	83.0	88.5	89.4	2015  CVPR
Mid-Level Filters 9	34.3	65.0	71.2	74.9	2014  CVPR
SalMatch 20	28.5	55.7	66.1	68.0	2014 ICCV

TOP r RANK MATCHING ACCURACY (%) ON CUHK01 DATASET.

# patch selection, feature fusion techniques is far more effective than just using saliency matching.

#### 512 4.2.4. Experiments on CUHK03

The CUHK03 is one of the highest cited person re-identification dataset 513 which consists of five different pairs of camera views, and the number of pictures 514 in this dataset exceeds 14,000. There are 13,164 bounding boxes detected by a 515 Deformable Part Model (DPM) of 1,467 different identities in CUHK03 dataset. 516 Fig. 7(d) and Table 5 provide the matching results of all the compared 517 algorithms. It can be seen that the proposed method is superior to the  $2^{nd}$ 518 best SPReID by 0.6% (91.8 - 91.2) in Rank-1, and ties with SPReID in Rank-519 20. SPReID extracts local features from human body parts obtained by hu-520 man semantic parsing. Both SPReID and our method use parsing for person 521 re-identification. SPReID focuses on parsing to make pedestrian body segmen-522 tation more accurate, while we focus on selecting reliable patch sequences for 523 precise matching. So if we learn from SPReID, the background noise will be 524 smaller. Then the selected patch sequence will be theoretically more reliable, 525

# Table 5:

Method	rank=1	rank=10	rank=15	rank=20	Reference
Ours	91.8	99.1	99.4	99.6	Proposed
BraidNet-CS+SRL 56	88.2	98.7	99.2	99.5	2018  CVPR
JointRe-id 13	54.7	91.5	96.8	97.3	2015  CVPR
FPNN 36	20.7	68.7	80.1	83.1	2014  CVPR
SPReID 57	91.2	99.2	99.5	99.6	2018  CVPR
EBb [41]	91.7	-	98.7	99.0	2018  CVPR
Ensembles 55	62.1	94.3	97.2	97.8	2015  CVPR
KISSME 2	14.2	52.6	66.4	70.0	2012  CVPR
LOMO+XQDA 3	52.2	92.1	95.6	96.3	2015  CVPR

TOP r RANK MATCHING ACCURACY (%) ON CUHK03 DATASET.

# and finally the accuracy will be improved.

- 527 4.2.5. Experiments on PRID 450S
- <sup>528</sup> The PRID 450S dataset contains 450 pairs of single-shot pedestrian images,
- <sup>529</sup> which are captured from two adjacent cameras. It is another challenging dataset,
- 530 similar to the VIPeR dataset, for background interference, partial occlusion and viewpoint changes.

Table 6:

TOP r RANK MATCHING ACCURACY (%) ON PRID 450S DATASET.

Method	rank=1	rank=10	rank=15	rank=20	Reference
Ours	72.5	96.4	97.8	98.7	Proposed
FFN <u>14</u>	66.6	92.8	96.6	96.9	2016 WACV
Kernel X-CRC 53	68.8	95.9	97.3	98.4	2019 JVCIR
LSSCDK 52	60.5	88.6	92.2	93.6	2016 CVPR
DRML 58	56.4	82.2	88.9	90.2	2016 ICIP
X-KPLS 59	52.8	90.0	94.8	95.4	2017 ICPR
MED_VL 60	45.9	82.9	89.8	91.1	2016 AAAI
ECM 4	41.9	76.9	82.6	84.9	2015 WACV

We evaluated the proposed approach by comparing the state-of-the-art ap-532 proaches on the PRID 450S dataset. This evaluation was conducted using the 533 images of detected persons. It can be seen from Table 6 that our method is su-534 perior to all compared state-of-the-arts, surpassing the  $2^{nd}$  best Kernel X-CRC 535 by 3.7% (72.5 - 68.8) in Rank-1, 0.5% (96.4 - 95.9) in Rank-10, 0.5% (97.8 - 97.3) 536 in Rank-15 and 0.3% (98.7 - 98.4) in Rank-20. Compared to Kernel X-CRC, 537 our local features contain gradient, color, and extreme points, not just the color 538 model as Kernel X-CRC does. It indicates the superiority of diverse features. 539 Fig. 7(e) describes the matching results of all the compared algorithms on the 540 PRID 450S dataset. 541

## 542 4.2.6. Experiments on iLIDS-VID

The iLIDS-VID dataset involves 300 different pedestrians observed across two disjoint camera views in a public open space. It comprises 600 image sequences of 300 distinct individuals, with one pair of image sequences from two camera views for each person. Each image sequence has variable length ranging from 23 to 192 image frames, with an average of 73 frames. The iLIDS-VID dataset is very challenging due to clothing similarities among people, lighting and viewpoint variations across camera views, cluttered background and random occlusions. Fig. [7(f)] and Table 7 show the matching results of all the

#### Table 7:

TOP r RANK MATCHING ACCURACY (%) ON ILIDS-VID DATASET.

Method	rank=1	rank=10	rank=15	rank=20	Reference
Ours	86.2	98.5	99.4	99.6	Proposed
CSsA+CSE 61	85.4	98.8	99.2	99.5	2018 CVPR
TDL <u>51</u>	56.3	95.6	97.9	98.3	2016 CVPR
TAPR <u>62</u>	55.0	93.8	96.9	97.2	2016 ICIP
$SI^2DL$ 63	48.7	89.2	96.6	97.3	2016 IJCAI
DRML 58	43.1	72.7	80.0	82.0	2016 ICIP
FAST3D 64	28.4	66.7	75.2	78.1	2016 ICIP

<sup>551</sup> compared methods. We can see that our method is superior to all the state-<sup>552</sup> of-the-art methods, surpassing the  $2^{nd}$  best method by 0.8% (86.2 - 85.4) in <sup>553</sup> Rank-1, 0.2% (99.4 - 99.2) in Rank-15, and 0.1% (99.6 - 99.5) in Rank-20. And <sup>554</sup> it's only 0.3% (98.8 - 98.5) lower than CSsA+CSE. Moreover, the proposed <sup>555</sup> method far surpasses recent methods (TDL, TAPR, SI<sup>2</sup>DL and DRML) in all <sup>556</sup> results. These validate that a combination of techniques may be more effective <sup>557</sup> than just using a single technique for person re-identification.

# 558 4.3. Ablation analysis

To further illustrate the rationality of each step of our process, we conduct ablation experiments for our method on VIPeR dataset. We verify the roles of four key parts of our algorithm, including background noise removal, saliency detection, local features, and global features through experiments. We take turns to remove key part (C1-C4) and compare with the complete method (C5), as shown in Table 8. It can be seen that Rank-1 and Rank-10 results of all incomplete methods are inferior to the complete method, which implies the importance of the default part.

#### Table 8:

# Experimental results for different configurations on VIPeR datasets.

Config.	Background noise removal	Saliency detection	Local features	Global features	Rank-1	Rank-10
C1	×	$\checkmark$	$\checkmark$	$\checkmark$	50.26	90.54
C2	$\checkmark$	×	$\checkmark$	$\checkmark$	43.31	86.52
C3	×	×	×	$\checkmark$	40.00	80.51
C4	$\checkmark$	$\checkmark$	$\checkmark$	×	51.74	89.62
C5	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	56.83	92.03

566

Firstly, we remove the process of background noise removal (C1), in other words, there is no primary selection, S1 = S. As can be seen from Table 8 that due to background noise, the performance degrades 6.57% and 1.49% for the Rank-1 and Rank-10 accuracy, respectively. Then we investigate the role of saliency detection (C2). We eliminate secondary selection based on saliency detection, and select 25 patches from S1 according to the principle of proximity.

Specifically,  $P(\cdot)$  in Eq (17) is redefined as the coordinate of the patch rather 573 than saliency value. It can be seen that the accuracy drops sharply, degrading 574 12.15% and 6.86% for the Rank-1 and Rank-10, respectively. Next, we remove 575 the local features (C3), which means that only global feature LOMO is used and 576 there is no patch selection. The feature representation is the same as 3. Rank-577 1 and Rank-10 become 40.00% and 80.51%. It also proves the necessity of the 578 patch selection method we proposed. Finally, global feature (C4) is removed to 579 demonstrate its role in compensating for inherent errors in saliency detection. 580 As can be seen from Table 8 that without the assistance of LOMO, there is 581 s slight decrease in accuracy, with rank-1 and rank-10 dropping by 5.09% and 582 2.42% respectively. 583

#### 584 4.4. Comparison with the most relevant methods

In this paper, the three key points of the proposed approach are utilizing local descriptors with the global features, background noise removal, and saliency detection. There are three corresponding algorithms, including LOMO+XQDA, saliency learning, and super-pixel segmentation for person re-identification. The following will introduce their differences with our proposed method and the experimental results.

# Table 9:

PERSON RE-ID MATCHING RATES(%) AT DIFFERENT RANKS ON VIPER, CUHK01, AND CUHK03 DATASETS

		VIPeR					CUHK01				CUHK03	
Method	rank@1	10	15	20	1	10	15	20	1	10	15	20
Ours	56.8	92.0	96.2	97.2	69.2	92.8	96.1	97.8	68.2	95.2	97.8	98.4
LOMO+XQDA 3	40.0	80.5	88.3	91.0	61.8	86.5	91.5	93.7	52.2	92.1	95.4	96.3
Saliency learning 28	44.1	81.8	88.4	91.2	28.5	55.7	66.4	68.0	56.8	93.8	96.2	97.5
BackSub-reid 26	27.2	64.2	75.2	77.8	19.2	44.8	65.8	68.7	40.2	72.1	84.5	86.4

590

We compared the matching rate with LOMO features, saliency learning, and the method based on super-pixel segmentation for person re-identification on VIPeR, CUHK01, and CUHK03 datasets. Table 9 records the results of the experiments which indicates that the proposed method is always better than others at rank-1.

596 4.5. Parameter analysis of the proposed method

#### <sup>597</sup> 4.5.1. The threshold of overlapping rate in background noise removal

The proposed system achieves accurate salient person re-identification through background removal based on super-pixel segmentation. However, in this paper, we apply the pedestrian parsing via a DDN network to achieve the background removal. The experiments show that the proposed method has obvious advantages over the other methods.



Figure 8: The relationship between matching rate and the threshold of overlapping rate (%) on VIPeR dataset.

602

In this paper, we take the pedestrian parsing as an important method for removing the background. We parse the pedestrians with the DDN network which allows the background to be removed from the edges of a human. It is an important preprocessing for picking up the pedestrian patches.

In the process of removing the background noise, we set the threshold of overlapping rate  $c(P_{i,j})$  to 25%, which was empirically determined after many experiments. It directly determines whether the patch belongs to the pedestrian or the background. The experimental result for choosing the overlap rate threshold are shown in the Fig. 8. We compared the matching rate when selecting different thresholds ranging from 0.1 to 0.4 at rank-1 on the VIPeR dataset, which shows that 25% as the threshold is appropriate.

#### <sup>614</sup> 4.5.2. The number of selected patches



Figure 9: The relationships between matching rate and the number of selected patches on VIPeR, CUHK01, CUHK03 datasets.

The number of patches in *S* has a great impact on the matching rate and execution efficiency. If the number is too small, effective information will be missed, resulting in lower accuracy. Too many will increase the computation time and reduce the execution efficiency.

Fig. 9 describes the relationship between the number and the matching rate at rank-1 on the VIPeR, CUHK01, CUHK03, and iLIDS-VID datasets. As we can see from Fig. 9, the matching rate increases as the number of selected patches increases. However, after the number exceeds 25, the rate of growth
becomes very slow, while the cost of time is multiplied. Finally, we selected 25
patches, which provides a compromise between computation time and matching
accuracy.

#### <sup>626</sup> 5. Conclusions

In this paper, we proposed a new patch selection method based on parsing 627 and saliency detection for person Re-identification. We solve the problem of 628 feature drift and patch imbalance of local features, and effectively compensate 629 for the inherent errors caused by saliency detection by combining local features 630 with global features. It provides more ideas for solving related problems. In 631 addition, our method can effectively deal with the real scenario of insufficient 632 samples, which has a strong engineering application value. It's another highlight 633 of our work. 634

# 635 Acknowledgments

The research is supported by National Natural Science Foundation of China(No.61471110, 61733003), National Key R&D Program of China(No.2017YFC0805005, 2017YFB1301103), Advance Research Project(41412050202), Fundamental Research Funds for the Central Universities(N172608005) and Natural Science Foundation of Liaoning(No.20180520040).

#### 641 References

- [1] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an
  ensemble of localized features," in *Computer Vision ECCV 2008, European Conference on Computer Vision, Marseille, France, October 12-18, 2008, Proceedings, 2008, pp. 262–275.*
- [2] M. Kstinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof, "Large
  scale metric learning from equivalence constraints," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2288–2295.

- [3] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local
   maximal occurrence representation and metric learning," in *IEEE Conference* on Computer Vision and Pattern Recognition, 2015, pp. 2197–2206.
- [4] X. Liu, H. Wang, Y. Wu, J. Yang, and M. H. Yang, "An ensemble color model
  for human re-identification," in *Applications of Computer Vision*, 2015, pp.
  868–875.
- [5] B. Ma, Y. Su, and F. Jurie, "Covariance descriptor based on bio-inspired
  features for person re-identification and face verification," *Image & Vision Computing*, vol. 32, no. 6-7, pp. 379–390, 2014.
- [6] S. Pedagadi, J. Orwell, S. Velastin, and B. Boghossian, "Local fisher dis criminant analysis for pedestrian re-identification," in *Computer Vision and Pattern Recognition*, 2013, pp. 3318–3325.
- [7] X. Y. Jing, X. Zhu, F. Wu, R. Hu, X. You, Y. Wang, H. Feng, and J. Y.
  Yang, "Super-resolution person re-identification with semi-coupled low-rank
  discriminant dictionary learning," *IEEE Transactions on Image Processing*,
  vol. 26, no. 3, pp. 1363–1378, 2017.
- [8] Y. Yang, J. Yang, J. Yan, S. Liao, D. Yi, and S. Z. Li, "Salient color names
  for person re-identification,", vol. 8689, no. 9, pp. 536–551, 2014.
- [9] R. Zhao, W. Ouyang, and X. Wang, "Learning mid-level filters for person
  re-identification," in *Computer Vision and Pattern Recognition*, 2014, pp.
  144–151.
- [10] Y. Wang, R. Hu, C. Liang, C. Zhang, and Q. Leng, "Camera compensation using a feature projection matrix for person reidentification," *IEEE Transactions on Circuits & Systems for Video Technology*, vol. 24, no. 8, pp.
  1350–1361, 2014.
- [11] D. Li, X. Chen, Z. Zhang, and K. Huang, "Learning deep context-aware
  features over body and latent parts for person re-identification," pp. 7398–
  7407, 2017.

- [12] A. Bosch, A. Zisserman, and X. Munoz, "Representing shape with a spa-677
- tial pyramid kernel," in ACM International Conference on Image and Video 678
- Retrieval, 2007, pp. 401-408. 679

694

- [13] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning archi-680 tecture for person re-identification," in Computer Vision and Pattern Recog-681 nition, 2015, pp. 3908-3916. 682
- [14] S. Wu, Y. C. Chen, X. Li, A. C. Wu, J. J. You, and W. S. Zheng, "An 683 enhanced deep feature representation for person re-identification," in Appli-684 cations of Computer Vision, 2016, pp. 1–8. 685
- [15] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: A 686 deep quadruplet network for person re-identification," in IEEE Conference 687 on Computer Vision and Pattern Recognition, 2017, pp. 1320–1329. 688
- [16] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable 689 person re-identification: A benchmark," in IEEE International Conference 690 on Computer Vision, 2015. 691
- [17] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan 692 improve the person re-identification baseline in vitro," in 2017 IEEE Inter-693 national Conference on Computer Vision (ICCV), Oct 2017, pp. 3774–3782.
- [18] L. Wei, S. Zhang, W. Gao, and Q. Tian, "Person transfer gan to bridge 695
- domain gap for person re-identification," in 2018 IEEE/CVF Conference on 696 Computer Vision and Pattern Recognition, June 2018, pp. 79–88. 697
- [19] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani, "Per-698 son re-identification by symmetry-driven accumulation of local features," in 699 Computer Vision and Pattern Recognition, 2010, pp. 2360–2367. 700
- [20] R. Zhao, W. Ouyang, and X. Wang, "Person re-identification by salience 701 matching," in IEEE International Conference on Computer Vision, 2014, pp. 702 2528 - 2535.703

- [21] B. Schlkopf, J. Platt, and T. Hofmann, "Graph-based visual saliency," in International Conference on Neural Information Processing Systems, 2006,
   pp. 545–552.
- [22] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng, "Person reidentification by multi-channel parts-based cnn with improved triplet loss
  function," in *Computer Vision and Pattern Recognition*, 2016, pp. 1335–1344.
- [23] H. Shi, Y. Yang, X. Zhu, S. Liao, Z. Lei, W. Zheng, and S. Z. Li, "Embedding deep metric for person re-identification: A study against large variations," 2016.
- [24] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Deep metric learning for person reidentification," in *International Conference on Pattern Recognition*, 2014, pp.
  34–39.
- <sup>716</sup> [25] Y. C. Chang, C. K. Chiang, and S. H. Lai, "Single-shot person re<sup>717</sup> identification based on improved random-walk pedestrian segmentation," in
  <sup>718</sup> International Symposium on Intelligent Signal Processing and Communica<sup>719</sup> tions Systems, 2013, pp. 1–6.
- [26] C. V. Le, Q. N. Hong, T. T. Quang, and N. D. Trung, "Superpixel-based
   background removal for accuracy salience person re-identification," in *IEEE International Conference on Consumer Electronics-Asia*, 2017, pp. 1–4.
- [27] P. Luo, X. Wang, and X. Tang, "Pedestrian parsing via deep decompositional network," in *IEEE International Conference on Computer Vision*,
  2014, pp. 2648–2655.
- [28] R. Zhao, W. Ouyang, and X. Wang, "Person re-identification by saliency
   learning," *IEEE Transactions on Pattern Analysis & Machine Intelligence*,
- vol. 39, no. 2, pp. 356–370, 2017.
- [29] Y. Chen, Z. Huo, and C. Hua, "Multi-directional saliency metric learning
  for person re-identification," *Iet Computer Vision*, vol. 10, no. 7, pp. 623–633,
  2017.

- [30] I. Kviatkovsky, A. Adam, and E. Rivlin, "Color invariants for person reiden-
- <sup>733</sup> tification," IEEE Transactions on Pattern Analysis & Machine Intelligence,
- vol. 35, no. 7, pp. 1622–1634, 2013.
- [31] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *International Conference on International Conference on Machine Learning*, 2010, pp. 807–814.
- [32] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P. A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion." *Journal of Machine Learning Research*,
  vol. 11, no. 12, pp. 3371–3408, 2010.

# [33] A. Z. A. O Chapelle, "Semi-supervised classification by low density sepa ration," AISTATS, 2005.

- [34] S. Liao, G. Zhao, V. Kellokumpu, M. Pietikainen, and S. Z. Li, "Modeling
  pixel process with scale invariant local patterns for background subtraction
  in complex scenes," in *Computer Vision and Pattern Recognition*, 2010, pp.
  1301–1306.
- [35] T. Ojala and I. Harwood, "A comparative study of texture measures with
  classification based on feature distributions," *Pattern Recognition*, vol. 29,
  no. 1, pp. 51–59, 1996.
- [36] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing
   neural network for person re-identification," in *Computer Vision and Pattern Recognition*, 2014, pp. 152–159.
- <sup>754</sup> [37] M. Hirzer, C. Beleznai, P. M. Roth, and H. Bischof, "Person re<sup>755</sup> identification by descriptive and discriminative classification," in *Scandina-*<sup>756</sup> vian Conference on Image Analysis, 2011.
- [38] L. Wei, Z. Rui, and X. Wang, "Human reidentification with transferred
   metric learning," in Asian Conference on Computer Vision, 2012.

- [39] P. M. Roth, M. Hirzer, M. Kstinger, C. Beleznai, and H. Bischof, Maha lanobis Distance Learning for Person Re-identification, 2014.
- [40] T. Wang, S. Gong, X. Zhu, and S. Wang, "Person re-identification by
   discriminative selection in video ranking," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 38, no. 12, pp. 2501–2514, 2016.
- [41] M. Tian, S. Yi, H. Li, S. Li, X. Zhang, J. Shi, J. Yan, and X. Wang,
  "Eliminating background-bias for robust person re-identification," 2018.
- [42] L. An, Z. Qin, X. Chen, and S. Yang, "Multi-level common space learning
  for person re-identification," *IEEE Transactions on Circuits & Systems for Video Technology*, vol. PP, no. 99, pp. 1–1, 2017.
- [43] D. Chen, Z. Yuan, B. Chen, and N. Zheng, "Similarity learning with spatial constraints for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1268–1277.
- Y. Yang, S. Liao, Z. Lei, and S. Z. Li, "Large scale similarity learning
  using similar pairs for person verification," in *Thirtieth AAAI Conference on*Artificial Intelligence, 2016, pp. 3655–3661.
- [45] X. Liu, M. Song, D. Tao, X. Zhou, C. Chen, and J. Bu, "Semi-supervised
  coupled dictionary learning for person re-identification," in *IEEE Conference*on Computer Vision and Pattern Recognition, 2014, pp. 3550–3557.
- [46] L. Ma, X. Yang, and D. Tao, "Person re-identification over camera networks
  using multi-task distance metric learning." *IEEE Transactions on Image Pro- cessing*, vol. 23, no. 8, pp. 3656–3670, 2014.
- <sup>781</sup> [47] Q. N. Hong, N. N. Tuan, T. T. Quang, D. N. Tien, and C. V. Le, "Deep
- rs2 spatio-temporal network for accurate person re-identification," in *Interna*-
- tional Conference on Information and Communications, 2017, pp. 208–213.
- [48] Y. Yan, B. Ni, Z. Song, C. Ma, Y. Yan, and X. Yang, "Person reidentification via recurrent feature aggregation," pp. 701–716, 2017.

- [49] J. Zhou, B. Su, and Y. Wu, "Easy identification from better constraints:
  Multi-shot person re-identification from reference constraints," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [50] S. Karanam, Y. Li, and R. J. Radke, "Person re-identification with dis criminatively trained viewpoint invariant dictionaries," in *IEEE International Conference on Computer Vision*, 2015, pp. 4516–4524.
- <sup>793</sup> [51] J. You, A. Wu, X. Li, and W. S. Zheng, "Top-push video-based person
  <sup>794</sup> re-identification," in *Computer Vision and Pattern Recognition*, 2016, pp.
  <sup>795</sup> 1345–1353.
- [52] Y. Zhang, B. Li, H. Lu, A. Irie, and R. Xiang, "Sample-specific svm learning
   for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1278–1287.
- [53] R. Prates and W. R. Schwartz, "Kernel cross-view collaborative representation based classification for person re-identification," 2016.
- [54] S. Li, M. Shao, and Y. Fu, "Cross-view projective dictionary learning for
   person re-identification," in *International Conference on Artificial Intelli- gence*, 2015, pp. 2155–2161.
- [55] S. Paisitkriangkrai, C. Shen, and A. V. D. Hengel, "Learning to rank in person re-identification with metric ensembles," in *Computer Vision and Pattern Recognition*, 2015, pp. 1846–1855.
- <sup>807</sup> [56] Y. Wang, Z. Chen, F. Wu, and G. Wang, "Person re-identification with cas<sup>808</sup> caded pairwise convolutions," in 2018 IEEE/CVF Conference on Computer
  <sup>809</sup> Vision and Pattern Recognition, June 2018, pp. 1470–1478.
- <sup>810</sup> [57] M. M. Kalayeh, E. Basaran, M. Gokmen, M. E. Kamasak, and M. Shah, <sup>811</sup> "Human semantic parsing for person re-identification," 2018.

- 812 [58] W. Yao, Z. Weng, and Y. Zhu, "Diversity regularized metric learning for
- person re-identification," in *IEEE International Conference on Image Pro*cessing, 2016, pp. 4264–4268.
- <sup>815</sup> [59] R. F. Prates and W. R. Schwartz, "Kernel hierarchical pca for person reidentification," in *International Conference on Pattern Recognition*, 2017.
- [60] A. F. O. A. Intelligence, "Association for the advancement of artificial
  intelligence," *Hyperfine Interactions*, vol. 6, no. 1, 2011.
- [61] D. Chen, H. Li, T. Xiao, S. Yi, and X. Wang, "Video person re-identification
  with competitive snippet-similarity aggregation and co-attentive snippet embedding," 2018.
- [62] C. Gao, J. Wang, L. Liu, J. G. Yu, and N. Sang, "Temporally aligned
  pooling representation for video-based person re-identification," in *IEEE International Conference on Image Processing*, 2016, pp. 4284–4288.
- [63] X. Zhu, X. Y. Jing, F. Wu, and H. Feng, "Video-based person reidentification by simultaneously learning intra-video and inter-video distance
  metrics," in *International Joint Conference on Artificial Intelligence*, 2016,
  pp. 3552–3558.
- [64] Z. Liu, J. Chen, and Y. Wang, "A fast adaptive spatio-temporal 3d feature
  for video-based person re-identification," in *IEEE International Conference*on Image Processing, 2016, pp. 4294–4298.

# 832 Authors



Yixiu Liu received the B.E. degrees from Northeastern University at Qinhuangdao, Qinhuangdao, China, in 2016, and is currently working toward the Ph.D. degree at School of School of Information Science and Engineering, Northeastern University, Shenyang, China. His research interests are in the area of computer vision, including person re-identification, pedestrian tracking.



838

839

846

849

850

851

Yunzhou Zhang received B.S. degree and M.S. degree in Mechanical and Electronic engineering from National University of Defense Technology, Changsha, China in 1997 and 2000, respectively. He received Ph.D. degree in pattern recognition and intelligent system from Northeastern University, Shenyang, China, in 2009, where he is currently a professor with at the School of Information Science and Engineering,

Northeastern University, China. His research interests include intelligent robot, 847 computer vision, and sensor networks. 848



in mathematics, statistics, and computing and the Ph.D. degree in mathematics from the University of Ulster, Londonberry, U.K., in 1999 and 2003, respectively. She is currently a Lecturer in the School of Computing and Intelligent System, Magee College, University of Ulster. She has more than 50 publications primarily in the field of mathematical image pro-

Sonya Coleman (M11) received the B.Sc. (Hons) degree

cessing, and much of the recent research undertaken by her has been supported 856 by funding from EPSRC award EP/C006283/11, the Leverhulme Trust, and the 857 Nuffield Foundation. Additionally, she is co-investigator on the EU FP7 funded 858 project RUBICON. She is the author or coauthor of over 70 research papers on 859 image processing, robotics, and computational neuroscience. Dr. Coleman was 860 awarded the Distinguished Research Fellowship by the University of Ulster in 861 recognition of her contribution to research in 2009. 862



Bir Bhanu (M82F95LF17) received the S.M. and E.E. degrees in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, MA, USA, the Ph.D. degree in electrical engineering from the University of Southern California, Los Angeles, CA, and the M.B.A. degree from the University of California at Irvine,

863

41

Irvine, CA. He was the Founding Professor of electrical en-869 gineering with the University of California at Riverside (UCR), Riverside, CA, 870 and served as its first Chair from 1991 to 1994. He has been the Cooperative 871 Professor of computer science and engineering (since 1991), bioengineering (s-872 ince 2006), and mechanical engineering (since 2008). He served as the Interim 873 Chair of the Department of Bioengineering from 2014 to 2016. He also served as 874 the Director of the National Science Foundation Graduate Research and Train-875 ing Program in video bioinformatics with UCR. He is currently the Bourns 876 Presidential Chair in engineering, the Distinguished Professor of electrical and 877 computer engineering, and the Founding Director of the Interdisciplinary Center 878 for Research in Intelligent Systems and the Visualization and Intelligent Sys-879 tems Laboratory, UCR. He has published extensively and has 18 patents. His 880 research interests include computer vision, pattern recognition and data min-881 ing, machine learning, artificial intelligence, image processing, image and video 882 database, graphics and visualization, robotics, human-computer interactions, 883 and biological, medical, military, and intelligence applications. In 1991, he was 884 a Senior Honeywell Fellow with Honeywell Inc. He is a Fellow of AAAS, IAPR, 885 SPIE, and AIMBE. 886



Shuangwei Liu received B.S. degree in automation from Northeastern University, Shenyang, China, in 2017. He is a graduate student in pattern recognition and intelligent systems in Northeastern University. He majors in computer vision and image processing, especially deep learning methods and person re-ID.

Yixiu Liu received the B.E. degrees from Northeastern University at Qinhuangdao, Qinhuangdao, China, in 2016, and is currently working toward the Ph.D. degree at School of School of Information Science and Engineering, Northeastern University, Shenyang, China. His research interests are in the area of computer vision, including person re-identification, pedestrian tracking.
Yunzhou Zhang received B.S. degree and M.S. degree in Mechanical and Electronic engineering from National University of Defense Technology, Changsha, China in 1997 and 2000, respectively. He received Ph.D. degree in pattern recognition and intelligent system from Northeastern University, Shenyang, China, in 2009, where he is currently a professor with at the School of Information Science and Engineering, Northeastern University, China. His research interests include intelligent robot, computer vision, and sensor networks.
Sonya Coleman received the B.Sc. (Hons) degree in athematics, statistics, and computing and the Ph.D. degree in mathematics from the University of Ulster, Londonberry, U.K., in 1999 and 2003, respectively. She is currently a Lecturer in the School of Computing and Intelligent System, Magee College, University of Ulster. She has more than 50 publications primarily in the field of mathematical image processing, and much of the recent research undertaken by her has been supported by funding from EPSRC award EP/C006283/11, the Leverhulme Trust, and the Nuffield Foundation. Additionally, she is co-investigator on the EU FP7 funded project RUBICON. She is the author or coauthor of over 70 research papers on image processing, robotics, and computational neuroscience. Dr. Coleman was awarded the Distinguished Research Fellowship by the University of Ulster in recognition of her contribution to research in 2009.
Bir Bhanu received the S.M. and E.E. degrees in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, MA, USA, the Ph.D. degree in electrical engineering from the University of Southern California, Los Angeles, CA, and the M.B.A. degree from the University of California at Irvine, Irvine, CA. He was the Founding Professor of electrical engineering with the University of California at Riverside (UCR), Riverside, CA, and served as its first Chair from 1991 to 1994. He has been the Cooperative Professor of computer science and engineering (since 1991), bioengineering (since 2006), and mechanical engineering (since 2008). He served as the Interim Chair of the Department of Bioengineering from 2014 to 2016. He also served as the Director of the National Science Foundation Graduate Research and Training Program in video bioinformatics with UCR. He is currently the Bourns Presidential Chair in engineering, the Distinguished Professor of electrical and computer engineering, and the Visualization and Intelligent Systems and the Visualization and Intelligent Systems Laboratory, UCR. He has published extensively and has 18 patents. His research interests include computer vision, pattern recognition and data mining, machine learning, artificial intelligence, image processing, image and video database, graphics and visualization, robotics,

human-computer interactions, and biological, medical, military, and intelligence
applications. In 1991, he was a Senior Honeywell Fellow with Honeywell Inc. He is a
Fellow of AAAS, IAPR, SPIE, and AIMBE.
Shuangwei Liu received B.S. degree in automation from Northeastern University, Shenyang, China, in 2017. He is a graduate student in pattern recognition and intelligent systems in Northeastern University. He majors in computer vision and image processing, especially deep learning methods and person re-ID.

Yixiu Liu
Yunzhou Zhang
Sonya Coleman
Bir Bhanu

Shuangwei Liu

\*Source Files - Latex or Word Click here to download Source Files - Latex or Word: Source files.zip We declare that we have no financial and personal relationships with other people or organizations that can inappropriately influence our work. There is no professional or other personal interest of any nature or kind in any product, service and/or company that could be construed as influencing the position presented in, or the review of, the manuscript entitled.