

# A Feature Extraction Method for Person Re-identification Based on a Two-branch CNN

Linlin Li

Renmin University of China

Bo Yang (✉ [13910700045@139.com](mailto:13910700045@139.com))

School of Emergency technology and management, North China Institute of Science & Technology,  
Langfang 065201, China <https://orcid.org/0000-0002-0961-3861>

Shaohui Chen

Beijing Gaocheng technology development co.LTD

---

## Research

**Keywords:** person re-identification, two-branch convolutional network, triplet loss function

**Posted Date:** October 7th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-84744/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# A feature extraction method for person re-identification based on a two-branch CNN

Linlin Li<sup>1</sup>, Bo Yang<sup>2\*</sup>, Shaohui Chen<sup>3</sup>

<sup>1</sup> School of information resources management, Renmin University of China, Beijing, 100872, China; xinxi\_li@tpri.org.cn (L.L.)

<sup>2</sup> School of Emergency technology and management, North China Institute of Science & Technology, Langfang 065201, China; 13910700045@139.com (B.Y.);

<sup>3</sup> Beijing Gaocheng technology development co. LTD, Beijing 100043, China; chensh01@ehualu.com (S.C);

\* Correspondence: 13910700045@139.com(B.Y.); Tel.: +8613910700045;

**Abstract:** A two-branch convolutional neural network (CNN) architecture for feature extraction in person re-identification (re-ID) based on video surveillance is proposed. Highly discriminative person features are obtained by extracting both global and local features. Moreover, an adaptive triplet loss function based on the original triplet loss function is proposed and is used in the network training process, resulting in a significantly improved learning efficiency. The experimental results on open datasets demonstrate the effectiveness of the proposed method.

**Keywords:** person re-identification; two-branch convolutional network; triplet loss function

## 1. Introduction

Recent years have witnessed a rapid development of the global economy and continuous progress of society, and thus, higher requirements have been put forward in various fields for video surveillance systems, the widespread use of which has brought much convenience for people's lives and has also made important contributions to public safety [1].

The targets with the highest frequency of occurrence in surveillance video are persons. Person re-identification (re-ID) refers to the analysis and judgement of the trajectory and range of motions of the targets of interest (such as criminal suspects and terrorists) by retrieving the obtained images of these targets in videos captured by other surveillance cameras in a timely manner to provide technical support and decision-making assistance to the government or the security sector. Common person re-ID tasks are generally composed of a query library and a search library. The query library contains the targets of interest. The search library consists of person images, which are generally obtained from videos using target detection algorithms. Specifically, the image in the query library is compared with each image in the search library, and the person image with the highest similarity is returned as the final recognition result. In the era of video data explosion, it is impossible to meet practical needs by relying only on manual data processing methods. Therefore, it is of important theoretical and practical significance to research and develop the corresponding person re-ID technology.

As early as in 1996, person re-ID began to draw the attention of researchers [2]. In 2006, Gheissari et al. [3] proposed for the first time the concept of person re-ID in an academic conference, followed by a surge of relevant studies. The VIPeR dataset is the first dataset specifically designed for person re-ID research [4] and greatly promoted the development of person re-ID field. The first related monograph authored by Gong et al. [5] was published in 2013, which discusses in detail the cutting-edge technologies and major challenges in the field of person re-ID. In recent years, researchers have begun to use deep learning to solve problems related to the person re-ID and have made great breakthroughs [6-9]. Many research outcomes have been published in major computer vision conferences and journals, and the recognition results on multiple datasets have also been significantly improved.

As the topic of person re-ID has attracted increasing attention, many solutions have been proposed, and good results have been achieved on different test datasets. However, person re-ID still

faces a number of challenges stemming from, for example, illumination changes, different person postures, varying viewing angles, non-aligned person images, complex image background, and scale changes.

Therefore, person re-ID in real situations is still one of hot topics in surveillance video research. In this study, we propose a two-branch convolutional neural network (CNN) architecture to extract global and local features of persons from surveillance videos to obtain person features with strong discriminative ability. Moreover, to address the problems in the network learning process, we propose an adaptive triplet loss function based on the conventional triplet loss function to greatly improve the learning efficiency. The experimental results on multiple open datasets verify the effectiveness of the proposed algorithm.

## 2. Related Works

One of the key issues in the person re-ID method based on feature extraction is to design a robust, reliable representation of person features, which should be able to not only identify different persons but also overcome the influence of complex environment on the identification. Current research methods in this category are mainly divided into manual feature extraction and deep learning-based feature extraction.

### 2.1. Manual feature extraction methods

In early research on person re-ID, shallow visual features, including color, shape, and trajectory features, were mainly used to solve the person re-ID problem. In 2012, Farenzena et al. proposed the symmetry-driven accumulation of local features (SDALF) [10]. According to the physiological structure of the human body, they divided a human into different parts, from which various color histogram features based on HSV color space were extracted and then combined into a whole for person matching. Later, Mignon et al. horizontally divided a person image into several blocks [11], extracted for each block the color features based on the YUV, HSV, and RGB color space in addition to the LBP texture features, and finally fused these different features into a whole to describe a specific person.

With the deepening of research, it was found that the use of shallow visual features alone was unable to well solve the person re-ID problem. Shallow visual features can partially represent the exterior intrinsic attributes of a person but are poorly adaptable to viewing angles and illumination. Therefore, researchers began to explore more effective representations of visual features, mainly including the semantic attributes and advanced visual features of a person.

In ECCV 2014, Yang et al. proposed the salient color names-based color descriptor (SCNCD) [12]. They argued that the clothing color of a person is very crucial for recognition; after extracting a variety of basic colors from a person, they extracted the corresponding color histograms from different areas of the person image and fused these color histograms as the final feature description. Liao et al. proposed in 2015 a feature descriptor named LOMO (local maximal occurrence) [13], described as follows. First, a person image is divided into six horizontal long stripes, and then a window of a certain size is used to move across each horizontal stripe to extract the HSV color histogram and SILTP histogram; the feature with the maximum value among these features is taken as the feature of a horizontal stripe, and finally the features of all horizontal stripes are combined as the LOMO feature descriptor. The LOMO feature is strongly invariant with respect to angle and illumination, making it widely applicable and often used for comparison with other algorithms. Recently, Matsukawa et al. proposed a feature descriptor named GOG (Gaussian of Gaussian) [14], in which the local area of an image is modeled by Gaussian distribution to simulate the appearance information of different local areas; very good results were obtained on multiple different datasets. In addition, some researchers also studied the use of semantic information and advanced visual features to represent a person. Because of the stability of the two, good recognition results can be achieved even if the posture of a person has changed considerably.

Although the abovementioned manual feature extraction methods lead to good results, they are mostly designed for certain specific situations and unable to achieve satisfactory test results in other

scenarios. In other words, manually extracted features have weak robustness and poor universality. Moreover, it is difficult to define the validity or applicable situation of each feature. In addition, most of the features are obtained using the multifeature fusion method, for which the fusion strategy of these features cannot ensure that the fused feature is optimal. Therefore, it is especially important to design more effective feature fusion strategies.

## 2.2. Feature extraction methods based on deep learning

The person feature extraction methods based on deep learning mainly use CNNs to extract person features. Compared to those extracted via the traditional manual feature extraction methods, the features extracted by the CNN model are relatively expressive, and thus the performance of recognition algorithms established with the CNN model will be substantially improved.

A common practice is to use the loss function as a constraint to train the parameters of the model to achieve the goal of “small intraclass distance and large interclass distance.” In 2016, Geng et al. proposed combining the classification loss and verification loss to train a network [15]. The main network has a two-stream CNN architecture, which is connected to the classification subnet and verification subnet. The classification subnet is used to predict the identity of the image, and the classification error loss is calculated based on the prediction results. The verification subnet fuses the features of the two images to determine whether the two images belong to the same person. During the test, the trained network is directly used to extract person features for re-ID. In 2017, Lin et al. noted that person identity information alone is not sufficient to learn a model with a high generalization ability. Therefore, they introduced person attribute labels through the labeled attribute information of a person so that the model needs to predict not only the person identity but also each person attribute correctly; the combined constraints of multiple features not only enhance the generalization ability of the model but also effectively improve its recognition results [16].

The rapid development of deep learning has promoted significant improvement in the person re-ID performance. At present, research on target feature extraction based on deep networks has the following deficiencies. First, the training datasets for person re-ID are generally small, and thus the trained network model tends to be overfitted, leading to insufficient generalization ability of person re-ID in real surveillance scenarios. Second, the deep features extracted based on deep learning networks are unable to effectively distinguish fine-grained target recognition; therefore, it is necessary to construct a new type of network to extract more essential target features.

## 3. Method

Due to the complexity and diversity of real surveillance scenarios, the effectiveness of traditional person re-ID methods based on manual features is far from satisfactory. Therefore, an increasing number of studies have focused on personal re-ID based on CNN. In 1989, LeCun et al. proposed for the first time a network capable of multilayer training named LeNet network [17], which was subsequently thoroughly studied by many researchers. ResNet is a more frequently used CNN for person re-ID.

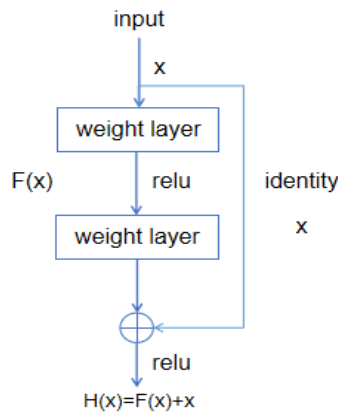
### 3.1 Residual neural network (ResNet) architecture

Researchers have already realized that the depth of the network is critical to the performance of the model. As the number of network layers is increased, the network can be expressed more effectively. Therefore, theoretically better results can be obtained with a deeper model. However, experimental findings reveal that a deep network is prone to the degradation problem, that is, as the network depth increases, the accuracy of the network becomes saturated or even decreases. This is because when the depth of CNN exceeds a certain number of layers, it will be very difficult to train the CNN due to the problems of gradient disappearance and gradient explosion, thus affecting the final recognition accuracy.

ResNet was proposed by He et al. of Microsoft Research in 2016 [18]. By using the residual unit, ResNet successfully trained a CNN with a depth of 152 layers, earning it the title of champion of the

ILSVRC 2015 competition, with a top-5 error rate of 3.75%. ResNet has a small number of parameters, and the residual unit designed by it can very quickly accelerate the training of neural networks. Moreover, the readiness of the model is also greatly improved.

For a common network with an input of  $x$  and an output of  $H(x)$ , the learning objective of the network is  $F(x) = H(x) - x$ . However, ResNet learns the difference between the input and output of the network, i.e., the residual  $H(x) - x$ . Such a residual hopping structure enables some information at the front end of the network to be directly transmitted to the back end of the network without the need for calculation in the middle layer. Therefore, the problem of gradient disappearance can be avoided during training so that the network can be trained very deeply. Figure 1 shows the schematic diagram of the residual unit in ResNet.



**Figure 1.** Schematic diagram of a residual unit

In essence, the residual unit is composed of several basic convolutional layers, but an identity connection is added between the input and output so that the information of different layers in the network can be fully utilized. ResNet is composed of such basic residual units, and the number of residual units used in the network represents the depth of the network. ResNet is referenced to the VGG19 network, and on this basis, downsampling is performed using convolution with a step size of 2, the full connection layer is replaced by the pooling layer, and the residual units are added through a short-circuit mechanism. These operations not only significantly reduce the number of parameters in ResNet but also improve the expressiveness of the network. Commonly used residual networks include ResNet18, ResNet34, ResNet50, ResNet101, and ResNet152.

ResNet50 is used as the backbone network in most of the deep learning-based person re-ID algorithms. To facilitate comparison with relevant algorithms, we also used ResNet50 as a backbone network to carry out relevant research.

### 3.2 Person feature extraction based on a two-branch CNN

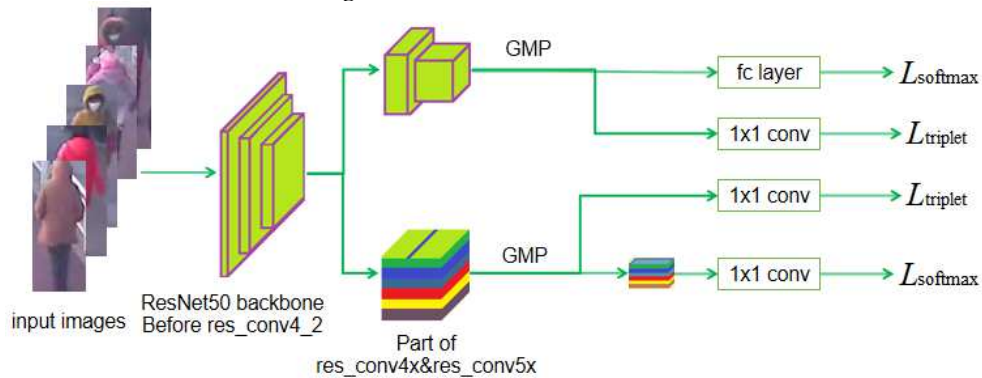
Good results were achieved via the use of CNNs for global person feature extraction during the early stage of person re-ID. With the development of research, it was found that use of global features only is not sufficient for person re-ID. Global features are effective for identifying persons with large differences in shape and color and poor otherwise. Therefore, researchers gradually paid attention to person re-ID based on local features. Attempts were made to extract local features of a person and then combine them with global features for person re-ID. Such a design can often achieve better results for person re-ID compared with use of only global features.

Currently, there are many methods for person re-ID based on combinations of global and local features. These methods focus on how to extract effective local features. Typical algorithms include Spindle Net [19], PDC [20], PL-Net [21], and GLAD [22]. Some of these algorithms directly divide a person image horizontally, extract the features separately from the divided image blocks, and finally combine them as a local feature. Using the key points of the human body for estimation, some other

algorithms first obtain the key points of persons in the image. Then, the image is divided into blocks for different parts of the human body based on the key points, features of different blocks of the human body image are extracted, and the features of these different parts of the body are combined to yield the final representation of the features. There are also algorithms that, according to the physiological structure of the human body and combining with the algorithm for estimation of human body key points, divide the person image into horizontal blocks for different parts, separately extract features from the divided image blocks, and finally fuse them into a local feature.

The above methods for the extraction of local features from a person image are simple, straightforward, easy to understand, and consistent with the human recognition process, thus also achieving good recognition results. However, it can be observed that these methods rely on the algorithm for estimation of human body key points, which is very time consuming and hence very inefficient for both the training process and practical applications. Moreover, in the process of feature extraction using the above algorithms, the local features and global features do not constrain and promote each other, and hence the learning efficiency must be improved.

Based on the above considerations, a CNN that can simultaneously extract the global and local features of a person is designed in this section. ResNet50 is used as the backbone network, and the first three layers of the network are used to extract the basic features of the image. Two branches are designed on the high-level semantic-level features to extract the global features and local features, respectively, of a person. The two parts work both collaboratively and separately, with weight sharing for the first three layers and independent weights for the subsequent advanced layers. In this manner, like the principle of human cognition of things, not only the overall information of a person can be seen but also the local information of different scales can be taken into account. Details of the network architecture are shown in Figure 2.



**Figure 2.** Schematic diagram of the feature extraction network architecture

In Figure 2, GMP stands for the global maximal pooling; 1x1 conv represents connection with a convolutional layer of size 1x1, which is used for feature dimension reduction;  $L_{triplet}$  denotes a triplet loss function; and  $L_{softmax}$  is a cross entropy loss function.

For the local feature extraction branch, the specific approach is as follows. After an image passes the first three layers of the ResNet50, the feature map is directly divided horizontally, and then each divided stripe area passes through the following convolutional layer, pooling layer, and 1x1 conv. The approach for horizontal division of the feature map simulates the division of the human body structure. Each stripe area of the feature map represents a different part of the human body. The full learning of these stripe areas is actually learning of the individualized areas of different parts of the human body, which is conducive to the learning of a discriminative feature of a person. Moreover, it can be seen that this design is in line with the human approach of judging a person. This is because in real life, it is often only necessary to determine who a person is based on parts of him or her rather than the appearance information. This recognition strategy of inferring the whole from the local is simple and effective, and the local branches in the designed network play exactly such a role. From another perspective, if the information based on each small piece can be recognized correctly, then

the accuracy of recognition by combining all small pieces of information will naturally increase significantly.

### 3.3 Design of network loss functions

The loss function of a network is used to evaluate the degree of difference between the predicted value of the model and the real value. It is often used as the objective function of a neural network and represents the optimization direction of an algorithm. In person re-ID, commonly used loss functions include the cross entropy loss function, comparison loss function [23], triplet loss function [24], and optimized triplet loss function [25].

The cross entropy loss function describes the distance between two probability distributions, and the smaller the cross entropy, the closer the two are. This loss function can lead to relatively good results for coarse-grained target classes. For fine-grained recognition tasks such as person re-ID, it can neither guarantee that the intraclass distance of samples is small enough nor that the interclass distance of samples is large enough. Therefore, the learning outcomes are very limited. Consequently, it is necessary to include the constraints on the spatial distribution of samples in the design of the loss function. Only by including the spatial constraints can the same class of samples be clustered in the feature space while different classes of samples be far apart from each other, which is conducive to the subsequent recognition. The triplet loss function based on the distance constraint between positive and negative samples is often used.

The triplet loss function [24] is a commonly used loss function for tasks such as retrieval and fine-grained recognition. Unlike the cross entropy loss function introduced above, the triplet loss function constraints the pairwise distance between samples; through constant iterative learning, it can decrease the intraclass distance and increase the interclass distance of samples in the feature space, thereby distinguishing different classes of persons.

The triplet loss function can decrease the distance between pairs of positive samples and increase the distance between pairs of negative samples. Finally, the person images of the same label form clusters in the feature space to distinguish different persons. However, we also observed that the triplet loss randomly sampled three images from the training data. Although this approach was simple, most of the samples were sample pairs that were simple and easy to distinguish. Thus, there were problems such as low training efficiency and nonideal convergence results. Therefore, it has been found that use of harder samples to train the network can make full use of the complex distribution characteristics of the training data, thereby improving the generalization ability of the network. Therefore, there appears a triplet loss function with batch hard mining.

The triplet loss function with batch hard mining (TriHard Loss function) is an optimization of the above basic triplet loss function [25]. By optimizing the network input triplets, each triplet involved in each training is optimal, thus improving the training efficiency and convergence speed of the network.

The core idea of TriHard Loss is as follows. To form a training batch,  $P$  persons with IDs are randomly selected, and  $K$  different images are randomly selected for each person, that is, one batch contains  $P \times K$  images. Then, for every image  $a$  in the batch, we can select the hardest positive sample and the hardest negative sample, which together with  $a$  form a triplet. First, the image set with the same ID as  $a$  is defined as  $A$ , and the remaining sets (with different ID) are  $B$ . Then, the TriHard Loss is expressed as

$$L_{\text{trihard}} = \frac{1}{P \times K} \sum_{a \in \text{batch}} (\max_{p \in A} d_{a,p} - \min_{n \in B} d_{a,n} + m)_+ \quad (1)$$

where  $m$  is the manually set interval between pairs of positive and negative samples,  $d_{a,p}$  represents the distance between the sample  $a$  and its positive samples,  $d_{a,n}$  represents the distance between the sample  $a$  and its negative samples, and  $()_+$  represents the greater value of the calculation result and zero. TriHard Loss fully utilizes the idea of batch hard mining in the loss calculation. Specifically, it calculates the distance between  $a$  and each image in feature space and then

selects the positive sample  $p$  with the longest distance from  $a$  and the negative sample  $n$  with the smallest distance from  $a$  to calculate the loss. Thus, TriHard Loss usually leads to better results than the traditional triplet loss.

We find that the interval between the positive and negative sample pairs in the triplet loss function is always fixed during the training. However, this fixed interval setting has an inevitable defect. When the interval is relatively large, the network needs to learn a relatively large interval in the feature space at the beginning of training, which is very difficult for the network. When the interval is relatively small, the network can be easily trained at the beginning and the learning outcomes are good; however, with the continuous deepening of learning, the setting of small intervals can greatly reduce the learning effect of the network, and at this time, a relatively large interval should be set instead to increase the learning difficulty of the network so that the network can continue to learn.

Second, we also find that during each training, TriHard Loss only selects the hardest samples of a sample  $a$  to participate in the training while ignore other samples, resulting in low training efficiency and insufficient learning. For the sample  $a$ , although other samples are not the hardest, they are the second hardest. With the continuous iteration of training, these second-hardest samples will gradually become the hardest and thus participate in the training. Therefore, for the sample  $a$ , during each training, we should consider that not only the hardest samples need to be fully learned but also other second-hardest samples need to participate in learning. Only in this manner can the learning efficiency of the network be improved and the convergence speed be accelerated.

Considering the above two points, we propose an adaptive TriHard Loss based on the original TriHard Loss.

First, we let the interval between pairs of positive and negative samples increase with the increasing number of iterations. The dynamic interval is expressed as follows:

$$m_{adaptive} = k \times epoch + b \quad (2)$$

where the left side of the equation represents the dynamic interval,  $k$  represents the degree of increase in interval,  $epoch$  represents the number of times that all data are trained in a round, and  $b$  represents an initial interval. With this design, the interval is small at the beginning of the training, and as the training continues, the interval gradually increases. Thus, during the whole training process, the feature intervals can all be fully learned, and the training efficiency will also be greatly improved.

Second, for every sample in the batch, regardless of whether it is a positive or negative sample, we assign an appropriate weight according to the distance between the two, so that every sample in the batch participates in training each time, greatly accelerating the efficiency of network learning. We refer to the improved TriHard Loss as the adaptive TriHard Loss, which is expressed as follows:

$$Loss = \frac{1}{P \times K} \sum_{a \in batch} \left( \sum_{p \in A} w_p d_{a,p} - \sum_{n \in B} w_n d_{a,n} + m_{adaptive} \right)_+ \quad (3)$$

where  $A$  represents the set of samples in a batch belonging to the same class as that of the sample  $a$  and  $B$  represents the set of samples in a batch belonging to different classes from that of the sample  $a$ . To calculate the specific loss for the sample  $a$ , we calculate the sum of weighted distances between the positive and negative sample pairs, respectively, formed with the sample  $a$  and then combine with the adaptive interval between sample pairs to obtain the loss of the sample  $a$ . The corresponding loss is calculated for each sample in the batch to obtain the overall loss.  $w_p$  and  $w_n$  are the weights assigned to the positive and negative samples, respectively, of sample  $a$ . They are expressed as follows:

$$w_p = \frac{e^{d_{a,p}}}{\sum_{pi \in A} e^{d_{a,pi}}}, w_n = \frac{e^{-d_{a,n}}}{\sum_{ni \in B} e^{-d_{a,ni}}} \quad (4)$$

The adaptive TriHard Loss can not only dynamically adjust the interval between positive and negative sample pairs but also set the adaptive weight. Therefore, the efficiency and stability of



training can both be greatly improved, and the convergence speed will also be faster, effectively reducing the risk of overfitting.

## 4. Results and discussion

In this section, we first explore the outcomes and effects of different forms of loss function and verify the effectiveness of the loss function proposed in this paper. Second, for the branches of local feature extraction, we explore the influence of dividing the feature map into different number of stripes on recognition. Finally, we also test the effectiveness and performance of different branches in the network.

### 4.1 Experimental configuration and parameter settings

The system environment of this study includes an Ubuntu 16.04 operating system, Intel Core i7-7700K CPU, 8 GB memory, NVIDIA GeForce GTX 1080 graphics card, PyTorch 0.4 open source deep learning framework, CUDA 8.0, and cuDNN 5.0.

The parameters during model training are as follows: the size of the input image =  $384 \times 128$ ; use of ADAM optimization algorithm for training, with the exponential decay rates of the first and second moment estimation, beta1 and beta2, being 0.9 and 0.999, respectively, and epsilon =  $10^{-8}$ ; initial learning rate =  $2 \times 10^{-4}$ , epoch = 200, and delay of learning rate at the 140th and 180th epochs, respectively, at a rate of  $5 \times 10^{-4}$ ; batch\_size = 48, K = 4, and P = 12.

In terms of experimental data, we mainly use the Market1501 dataset, and the subsequent experimental environment and parameters are basically consistent with the above settings.

### 4.2 Evaluation indices for person re-ID algorithm

The commonly used indices for evaluating person re-ID algorithms include rank-k and mAP.

(1) rank-k

The rank-1 matching rate (rank-1) [26] is a commonly used index, which refers to the probability that the image to be queried and the image that ranks the first in similarity in the search library belong to the same target, i.e.,

$$\text{rank} - 1 = \frac{\sum_{i \in \{1, 2, \dots, m\}} S_i}{m} \quad (5)$$

where m is the total number of images and  $S_i$  is a flag variable representing whether the ith image to be queried and the image ranking the first in similarity belong to the same target ( $S_i = 1$  if so, and  $S_i = 0$  if not). In general, a larger rank-1 means a better performance of the model. Therefore, rank-1 is the most direct and the most important index, and it has been used extensively to evaluate the performance of the model.

The k-matching rate (rank-k) denotes the probability that the image in the top k-position of similarity ranking in the retrieval database belongs to the same pedestrian as the image to be retrieved. Commonly used are rank-1, rank-3, rank-5, rank-10, and rank-20. As with rank-1, a larger rank-k means a better performance of the model. The rank-k index represents the judgment of whether there exists at least one image in the first k images belonging to the same person as the image to be queried. Hence, rank-k is a reflection of the comprehensive search ability of the model and can more comprehensively measure the performance of the model than rank-1.

(2) mAP

Because the rank-k index cannot well measure the recall rate of the model, Zheng et al. in 2015 introduced for the first time the mean average precision (mAP) index into the evaluation system for person re-ID [27]. The mAP index is a trade-off between the precision rate and the recall rate and can more objectively and comprehensively evaluate the performance of the model.

Based on the combinations of real categories and model prediction categories, Common classification problems can be divided into four cases, namely, true positive (TP), false positive (FP), true negative (TN), and false negative (FN). The precision rate P and the recall rate R are as follows:

$$P = \frac{TP}{TP + FP} \quad (6)$$

$$R = \frac{TP}{TP + FN} \quad (7)$$

Then, a two-dimensional curve (called P-R curve) can be drawn by taking the precision rate P as the vertical coordinate and the recall rate R as the horizontal coordinate. The area enclosed under the P-R curve is called the average prevision (AP), and the mean of AP values of all classes is called mAP.

#### 4.3 Experimental comparison of loss functions

In this section, the effects of different loss functions are tested. ResNet50 is used as the backbone network. The outputs of the last convolutional layer of ResNet50 are used as the person features. Then, the Euclidean distance is used to calculate the similarity for recognition. The parameters of different loss functions remain consistent during the training. Market1501 [27] is used as test dataset. The test results of different loss functions are reported in Table 1.

**Table 1** Test results of different loss functions

Loss function	mAP (%)	rank1(%)
<b>Softmax</b>	41.3	65.8
<b>Triplet loss</b>	54.8	75.9
<b>TriHard Loss</b>	68.0	83.8
$L_{adaptive}$	69.7	85.2

It can be observed that because the triplet loss constraints the distance between sample pairs, it has better test results than the softmax loss, which also shows that it is not enough to learn different classes only; rather, an in-depth learning of the distance relationship between classes is also required. Second, comparison of the triplet loss and the TriHard Loss finds that both mAP and rank1 indices are greatly improved, indicating that batch hard mining can be of great hep to the triplet loss in that improving the quality of triplet sampling ensures the availability of high-quality triplets every time to train the network to improve the training results. Finally, a comparison of our proposed adaptive TriHard Loss and the original TriHard Loss reveals that mAP and rank-1 are increased by 1.7% and 1.4%, respectively, indicating the effectiveness of our proposed adaptive TriHard Loss.

#### 4.4 Experimental results of person feature extraction

As is known from the introduction of the network architecture in earlier sections, local feature extraction requires horizontal division of the feature map. Therefore, experiments are carried out on the number of horizontal stripes resulting from the division of the feature map to obtain the optimal design of the network architecture. The detailed experimental results are reported in Table 2.

**Table 2** Effect of the number of divided stripes in the feature map on the network performance

Number of local regions divided, d	mAP (%)	Rank-1(%)
1	75.4	88.5
2	79.7	90.3
4	81.8	92.7
6	83.4	93.2
8	81.6	92.4
12	80.5	90.8

It can be found from observation of the above experimental results that with increasing number of divided stripes, both the mAP and rank-1 values of the model first increase and then decrease, a variation pattern that is consistent with common sense. The smaller the number of divided stripes in the feature map, the closer the results to those of the global feature. As a coarse-grained person

descriptor, the global feature can only be used to compare the overall appearance and outline of different persons, but it is not sufficient to distinguish persons that require local fine comparison, resulting in a low recognition rate at the beginning. With increasing number of divided stripes, the model will pay more attention to the individualized area in each of the different stripes to learn more discriminative local features, and therefore the results will be increasingly better. However, as the number of divided stripes increases, the information contained in each stripe area decreases, leading to disordered network learning and lowered learning outcomes. Therefore,  $d = 6$  is generally used as the default setting.

To illustrate the effectiveness of the network designed in this paper, we conduct experiments on the global feature branch, local feature branch, and the combination of the two, respectively, in the network to clarify the function of each branch. The loss function is set as  $1 \times \text{CrossEntropy Loss} + 2 \times \text{Adaptive TriHard Loss}$ , where  $1 \times \text{CrossEntropy Loss}$  represents the use of one times the cross entropy loss function, and  $2 \times \text{Adaptive TriHard Loss}$  represents the use of two times the adaptive TriHard Loss. Hence, the final total loss function contains both the cross entropy loss function and the adaptive TriHard Loss, making the overall learning outcome of the network better. The detailed experimental results are reported in Table 3.

**Table 3** Test results of different branches of the network

Networks with different branches	mAP (%)	rank1(%)
Global_feature	75.4	88.5
local_feature	83.4	93.2
global_feature + local_feature	84.6	93.8

where global\_feature represents the global feature branch and local\_feature represents the local feature branch. First, it can be found that local features have a significant effect on the mAP and rank-1 indices, indicating that the local feature branch effectively learns the discriminative region of a person and verifies the effectiveness of the local feature branch. Second, it can be observed that the combination of the global feature branch and the local feature extraction branch achieves the best results, indicating that although the local feature branch is highly efficient, there is still some overall information that is not learned, and the global feature branch can exactly make up for this defect. Therefore, it is illustrated that the two are highly complementary and can promote each other to jointly improve the recognition outcomes.

Moreover, we compare the multibranch network proposed in this paper with a similar algorithm published recently (see Table 4 and Table 5) using the test datasets of CUHK03 [28] and Market1501 [27].

**Table 4** Test results of different algorithms on the CUHK03 dataset

Algorithm	Rank-1(%)	Rank-5(%)	Rank-10(%)
Spindle Net[19]	88.5	97.8	98.6
PDC[20]	88.7	98.6	99.2
PL-Net[21]	82.8	96.6	98.6
GLAD[22]	85.0	97.9	99.1
This paper	90.6	98.7	99.4

**Table 5** Test results of different algorithms on Market1501 dataset

Algorithm	mAP (%)	rank1(%)
Spindle Net[19]	---	76.9
PDC[20]	63.4	84.1

PL-Net[21]	69.3	88.2
GLAD[22]	73.9	89.9
This paper	84.6	93.8

As can be observed from the above tables, compared with other algorithms that are based on combinations of global and local features, the network proposed in this paper achieves the best results, and values of each index far exceed those of the other algorithms, indicating the effectiveness of the multibranch network.

## 5. Conclusion

This paper first introduces a CNN-based person feature extraction method; second, a multibranch network architecture combining a global feature branch and a local feature branch is constructed, and the construction of each branch of the network is described in detail; third, the loss functions commonly used in person re-ID are discussed, and an adaptive triplet loss function is proposed. In the experiments, the effects of various types of loss functions are validated, the rationality of the design of the multibranch network architecture is also demonstrated, and the effectiveness of the setting of global feature branches and local feature branches is verified.

**Funding:** This research was funded by the National Natural Science Foundation of China (Grant Nos. 41671441, 41531177, and U1764262). This research was funded by the National Key Research and Development Project (Grant Nos.2018YFB1600600)

**Acknowledgments:** We are grateful for the assistances of the reviewers and editors.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## Abbreviations

<b>CNN:</b>	Convolutional neural network
<b>re-ID:</b>	Re-identification
<b>SDALF:</b>	Symmetry-driven accumulation of local features
<b>LBP:</b>	Local Binary Pattern
<b>SIFT:</b>	Scale invariant feature transform
<b>PDC:</b>	Pose-driven Deep Convolutional
<b>PL-Net:</b>	Part loss net
<b>GLAD:</b>	Global-local-alignment descriptor
<b>SCNCD:</b>	Salient color names-based color descriptor
<b>LOMO:</b>	Local maximal occurrence
<b>GOG:</b>	Gaussian of Gaussian
<b>GMP:</b>	Global max pooling

## References

- Li Y, Wang X. Investigation of intelligent video surveillance system based on artificial intelligence technology [J]. Technology Innovation and Application, 2018(34): 76-77.
- Cai Q, Aggarwal J K. Tracking human motion using multiple cameras [C]. International Conference on Pattern Recognition. Vienna, Austria, 1996: 68-72.
- Gheissari N, Sebastian T B, Hartley R. Person re-identification using spatiotemporal appearance [C]. IEEE Conference on Computer Vision and Pattern Recognition. New York, USA, 2006: 1528-1535.
- Gray D, Brennan S, Tao H. Evaluating appearance models for recognition, reacquisition, and tracking [J]. International journal of computer vision, 2007, 89(2): 56-68.
- Gong S G, Cristani M, Yan S C, et al. Person re-identification [J]. Advances in Computer Vision and Pattern Recognition, 2013, 42(7): 301-313.

- 461 6. W.S. Zheng, X. Li, T. Xiang, et al. Partial person re-identification. in ICCV, 2015.
- 462 7. Yi D, Lei Z, Li S Z, Deep metric learning for practical person re-identification [C].International Conference  
463 on Pattern Recognition. Stockholm Waterfront, Sweden, 2014.
- 464 8. Oreifej O, Mehran R, Shah M. Human identity recognition in aerial images [C]. IEEE Conference on  
465 Computer Vision and Pattern Recognition. San Francisco, USA, 2010: 709-716.
- 466 9. Jungling K, Bodensteiner C, Arens M. Person re-identification in multi-camera networks [C]. Computer  
467 Vision and Pattern Recognition Workshops. Colorado, USA, 2010: 709-716.
- 468 10. Farenzena M, Bazzani L, Perina A, et al. Person re-identification by symmetry-driven accumulation of local  
469 features [C] // Computer Vision and Pattern Recognition(CVPR), 2010 IEEE Conference on. 2010: 2360-2367.
- 470 11. Mignon A, Jurie F. PCCA: A new approach for distance learning from sparse pairwise constraints [C]//2012  
471 IEEE conference on computer vision and pattern recognition. IEEE, 2012: 2666-2672.
- 472 12. Yang Y, Yang J, Yan J, et al. Salient color names for person re-identification [C]// European conference on  
473 computer vision. 2014: 536-551.
- 474 13. Liao S, Hu Y, Zhu X, et al. Person re-identification by local maximal occurrence representation and metric  
475 learning [C]// the IEEE Conference on Proceedings of Computer Vision and Pattern Recognition. 2015: 2197-  
476 2206.
- 477 14. Matsukawa T, Okabe T, Suzuki E, et al. Hierarchical gaussian descriptor for person re-  
478 identification[C]//IEEE Conference on Computer Vision and Proceedings of the Pattern Recognition. 2016:  
479 1363-1372.
- 480 15. Geng M, Wang Y, Xiang T, et al. Deep Transfer Learning for Person Re-identification [J]. 2016.
- 481 16. Lin Y, Zheng L, Zheng Z, et al. Improving Person Re-identification by Attribute and Identity Learning [J].  
482 2017.
- 483 17. LeCun Y, Boser B, Denker J S, et al. Backpropagation applied to handwritten zip code recognition [J]. Neural  
484 computation, 1989, 1(4): 541-551.
- 485 18. He K, Zhang X, Ren S, et al. Deep residual learning for image recognition [C] //Proceedings of the IEEE  
486 conference on computer vision and pattern recognition. 2016: 770-778.
- 487 19. Zhao H, Tian M, Sun S, et al. Spindle net: Person re-identification with human body region guided feature  
488 decomposition and fusion [C]//Proceedings of the IEEE Conference on Computer Vision and Pattern  
489 Recognition. 2017: 1077-1085.
- 490 20. Su C, Li J, Zhang S, et al. Pose-driven Deep Convolutional Model for Person Re-identification [C]//2017  
491 IEEE International Conference on Computer Vision (ICCV). 2017: 3980-3989.
- 492 21. Yao H, Zhang S, Hong R, et al. Deep representation learning with part loss for person re-identification [J].  
493 IEEE Transactions on Image Processing, 2019.
- 494 22. Wei L, Zhang S, Yao H, et al. Glad: Global-local-alignment descriptor for person retrieval [C]//Proceedings  
495 of the 25th ACM international conference on Multimedia. ACM, 2017: 420-428.
- 496 23. Varior R R, Shuai B, Lu J, et al. A siamese long short-term memory architecture for human re-identification  
497 [C]//European Conference on Computer Vision. Springer, Cham, 2016: 135-153.
- 498 24. Schroff F, Kalenichenko D, Philbin J. Facenet: A unified embedding for face recognition and clustering  
499 [C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2015: 815-823.
- 500 25. Hermans A, Beyer L, Leibe B. In Defense of the Triplet Loss for Person Re-Identification [J]. 2017.
- 501 26. Bolle R M, Connell J H, Pankanti S, et al. The relation between the ROC curve and the CMC [C]//Fourth  
502 IEEE Workshop on Automatic Identification Advanced Technologies (AutoID'05). IEEE, 2005: 15-20.
- 503 27. Chen Y C, Zheng W S, Lai J H, et al. An asymmetric distance model for cross-view feature mapping in  
504 person reidentification [J]. IEEE transactions on circuits and systems for video technology, 2017, 27(8): 1661-  
505 1675.

28. Li W,Zhao R,Xiao T,et al. DeepReID :Deep filter paring neural network for person re-identification [C].  
Proceeding of the IEEE Computer Society Conference on Computer Vision and Person Recognition.  
Columbus, Ohio:2014: 152-159.

#### **Availability of data and materials**

Please contact author for data requests.

#### **Competing interests**

The authors declare that they have no competing interests.

#### **Funding**

This research was funded by the National Natural Science Foundation of China (41671441, 41531177, U1764262). This research was funded by the National Key Research and Development Project (Grant Nos.2018YFB1600600)

#### **Authors` Contributions**

All authors took part in the discussion of the work described in this paper. All authors read and approved the final manuscript.

#### **Acknowledgements**

Thanks for the help of reviewers and editors.

#### **Authors` information**

# Figures

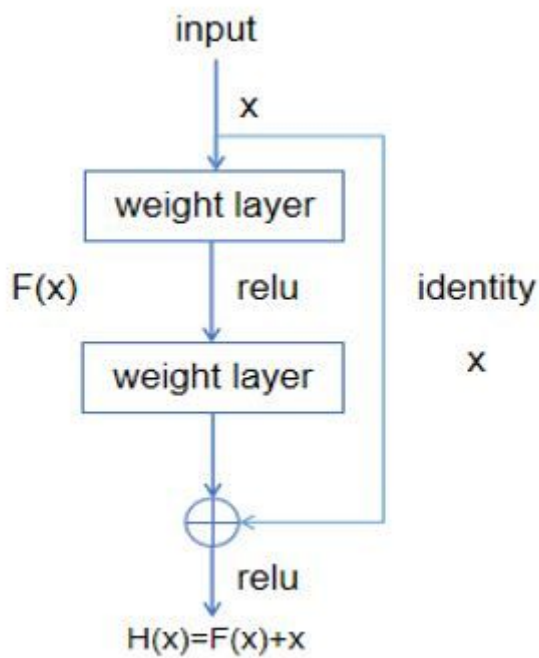


Figure 1

Schematic diagram of a residual unit

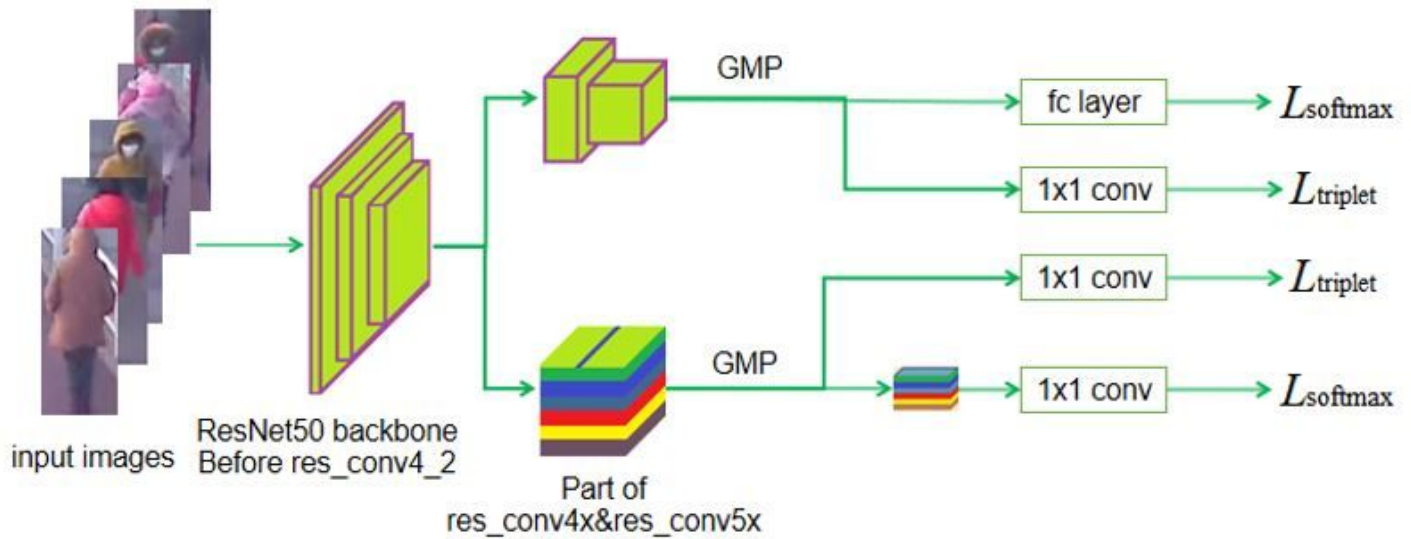


Figure 2

Schematic diagram of the feature extraction network architecture