

MHSA-Net: Multi-Head Self-Attention Network for Occluded Person Re-Identification

Hongchen Tan, Xiuping Liu, Baocai Yin and Xin Li, *Senior Member, IEEE*

Abstract—This paper presents a novel person re-identification model, named Multi-Head Self-Attention Network (MHSA-Net), to prune unimportant information and capture key local information from person images. MHSA-Net contains two main novel components: Multi-Head Self-Attention Branch (MHSAB) and Attention Competition Mechanism (ACM). The MHSAB adaptively captures key local person information, and then produces effective diversity embeddings of an image for the person matching. The ACM further helps filter out attention noise and non-key information. Through extensive ablation studies, we verified that the Multi-Head Self-Attention Branch (MHSAB) and Attention Competition Mechanism (ACM) both contribute to the performance improvement of the MHSA-Net. Our MHSA-Net achieves competitive performance in the standard and occluded person Re-ID tasks.

Index Terms—Occluded Person Re-ID, Multi-Head Self-Attention, Attention Competition Mechanism, Feature Fusion.

I. INTRODUCTION

Person re-identification (Re-ID) is a fundamental task in distributed multi-camera surveillance. It identifies the same person in different (non-overlapping) camera views. Re-ID has important applications in video surveillance and criminal investigation. With the surge of interest in deep representation learning, the person Re-ID task has achieved great progress in recent years [61]. Although recently many methods [52], [72], [71], [12], [77], [81], [58] have boosted the performance of the standard person Re-ID task, they didn't consider the situation that the person is occluded by various obstructions like cars, trees, or other people. The occlusion in person images is still a key challenging issue that hinders Re-ID performance. Thus, this paper aims to develop a Re-ID algorithm that can better handle occlusions in images.

In the occluded person Re-ID task, occluded regions often contain a lot of noise that results in mismatching. So a key issue in occluded Re-ID is to build discriminative features from unoccluded regions. Some part-based methods [55], [39], [41] manually crop the occluded target person in probe images and

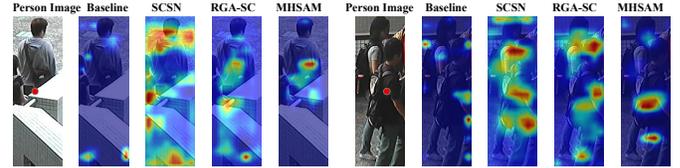


Fig. 1. The occluded person images' attention maps are produced by our Baseline, RGA-SC [66], SCSN (3-stage) [6] and the person Re-ID model equipped with Multi-Head Self-Attention mechanism (MHSAM) [3], [25]. The red dot is the target person.

then use the unoccluded parts as the new query. However, these manual operations are inefficient in practice. Another type of approach is to use human model to help build person features. More recently, [29], [18], [17] applied pose estimators to obtain the person's key points to locate effective regions of the person. However, the difference between training datasets of pose estimation and that of person retrieval often exist, making pose estimation based feature extraction sometimes unstable. It is desirable to design an effective mechanism to adaptively capture the key features from non-occlusion regions without relying on human models.

We are inspired by the recent Multi-Head Self-Attention mechanism (MHSAM) [3], [25], [34], [67], which flexibly captures spatially different local salience from the whole image, and generates multiple attention maps, from different aspects, for a single image. With MHSAM, noisy/unimportant regions can be pruned and key local feature information can be highlighted. Therefore, we believe the idea of MHSAM can help a Re-ID model to better locate key features from occluded images. As shown in Fig. 1, compared with the Baseline, two outstanding attention Re-ID model RGA-SC [66] and SCSN (3-stage) [6], the MHSAM can help the person Re-ID model better capture key information of the target person from the unoccluded regions and avoid information from occluded regions. The baseline may undesirably pay attention to clutter regions (left example) or other persons (right example), while our MHSAM model handles such occlusions much better.

However, developing effective MHSAM for the task of Re-ID is non-trivial and needs careful design. We propose a novel MHSAM module for the person Re-ID task with a set of new strategies to help select the key sub-regions in the image. We call our new attention module a *Multi-Head Self-Attention Branch* (MHSAB).

Furthermore, attention noise from the occluded or non-key regions often exists, and it affects the performance of the Re-ID model, we propose to design an *Attention Competition*

(Corresponding author: Baocai Yin)

Hongchen Tan and Baocai Yin are with Artificial Intelligence Research Institute, Beijing University of Technology, Beijing 100124, China (e-mail: tanhongchenphd@bjut.edu.cn; ybc@bjut.edu.cn).

Xin Li is with School of Electrical Engineering & Computer Science, and Center for Computation & Technology, Louisiana State University, Baton Rouge (LA) 70808, United States of America (e-mail: xinli@cct.lsu.edu).

Xiuping Liu is with School of Mathematical Sciences, Dalian University of Technology, Dalian 116024, China (xpliu@dlut.edu.cn).

Mechanism (ACM) to further help MHSAB suppress or filter out such attention noise from non-key sub-regions. Our main contributions are as follows:

- (I) We proposed a new attention module, MHSAB, that can more effectively extract person features in occluded person Re-ID.
- (II) We proposed a new attention competition module (ACM) to better prune attention noise from unimportant regions.
- (III) By integrating MHSAB and ACM modules, our final MHSA-Net framework demonstrates better performance over most state-of-the-art methods when processing occluded images, i.e., on four occlusion datasets: Occluded-DukeMTMC [29], P-DukeMTMC-reID [30], Partial-REID [55], and Partial-iLIDS [39]. On standard generic person Re-ID datasets, e.g., Market-1501 [38], DukeMTMC-reID [46], [74], and CUHK03 [53], our MHSA-Net produces similar results with these state-of-the-art algorithms.

II. RELATED WORK

A. Attention Mechanism in Person Re-identification

Attention mechanisms have been widely exploited in computer vision and natural language processing, for instance in Text-to-Image Synthesis [22], Object Tracking [2], Image/Video Captioning [28], Visual Question Answering [82], Neural Machine Translation [15], and some Video Tasks [69], [37], [68]. It can effectively capture task-relevant information and reduce interference from less important ones. Recently, many person Re-ID approaches [62], [33], [10], [60], [8], [32], [51], [5] also introduced various attention mechanisms into deep models to enhance identification performance.

[33], [10], [60], [8] applied a human part detector or a human parsing model to capture features of body parts. [14] explored both the human part masks and human poses to enhance human body feature extraction. [32], [29] exploited the connectivity of the key points to generate human part masks and focuses on the human's representation. However, the success of such approaches heavily relies on the accuracy of the human parsing models or pose estimators.

Other methods typically focus on extracting the person appearance or gait information, from the 3D space or depth images, to reduce the interference of background or occlusion. For example, Zhedong et al. [76] try to project the 2D person image into the 3D space, and conduct the person matching in the 3D space. Munaro et al. [42] proposed point cloud matching (PCM) strategy to compute the distances of multi-view point cloud sets, so as to distinguish between different persons. Haque et al. [20] adopted 3D LSTM to build motion dynamics of 3D person point clouds for person matching. [45] proposed a self-supervised gait encoding approach that can leverage unlabeled 3D skeleton data to learn gait representations for person Re-ID. Sivapalan et al. [50] extended the Gait Energy Image (GEI) [11] to 3D domain and proposed Gait Energy Volume (GEV) strategy based on depth images to perform gait-based person Re-ID. In [35], **Convolutional Neural Network**

Long Short-Term Memory (CNN-LSTM) with reinforced temporal attention (RTA) was proposed for person matching based on a split-rate RGB-Depth transfer method.

Besides, many methods [51], [5], [54], [9], [63] tried to exploit a different type of attention mechanism that does not need to use human models to capture human body features. [27] proposed a dual attention matching network based on an inter-class and an intra-class attention module to capture context information of video sequences for person Re-ID. ABD-Net [51] combined spatial and channel attention to directly learn human's information from the data and context. [63] calculated the similarity of the local features to enhance local part information. [9] applied an attribute classification to gain local attention information. However, it does not consider how to filter out information from the occlusion regions in the image. Therefore, with its fixed and parameter-free attention patterns, information from the occlusion region will be inevitably included.

Similar to [51], [27], [9], [63], our attention module also does not rely on an external human model. But different from these methods, our attention mechanism can adaptively enhance/suppress attention weights of local features through a multi-parameter learning strategy. The attention information of occluded and unoccluded regions in our attention mechanism is adaptively adjusted according to the targeting task. For the person Re-ID task, our attention module can flexibly capture the key local features and prune out information from occlusion regions.

B. Occluded Person Re-identification

Occlusion is a key challenging issue in person Re-ID. Recent studies [55], [39], [41], [30], [29], [40], [17], [18] on this topic can be divided into two categories: (i) partial person Re-ID methods [55], [39], [41], and (ii) occluded person Re-ID methods [30], [29], [40], [17], [18].

The former category aims to match a partial probe image to a gallery holistic image. For example, [55] adopted a global-to-local matching mechanism to capture the key information from the spatial channel of the feature maps. DSR [39], [41] proposed a spatial feature reconstruction strategy to align the partial person image with holistic images. However, these methods need a manual crop of the occluded target person in the probe image, before the cropped unoccluded part can be used to retrieve the target person.

The latter category aims to directly capture key features from the whole occluded person image to perform the person matching. AFPB [30] combined the occluded/unoccluded classification task and person ID classification task to improve the performance of deep model on capturing key information. FPR [40] reconstructed the feature map of unoccluded regions in occluded person image and further improved it by a foreground-background mask to avoid the influence of background clutter. [29], [17], [18] proposed pose guided feature alignment methods to match the local patches of query and the gallery images based on human key-points. Our MHSA-Net also belongs to this type of method. However, different from these methods, the MHSA-Net does not require any

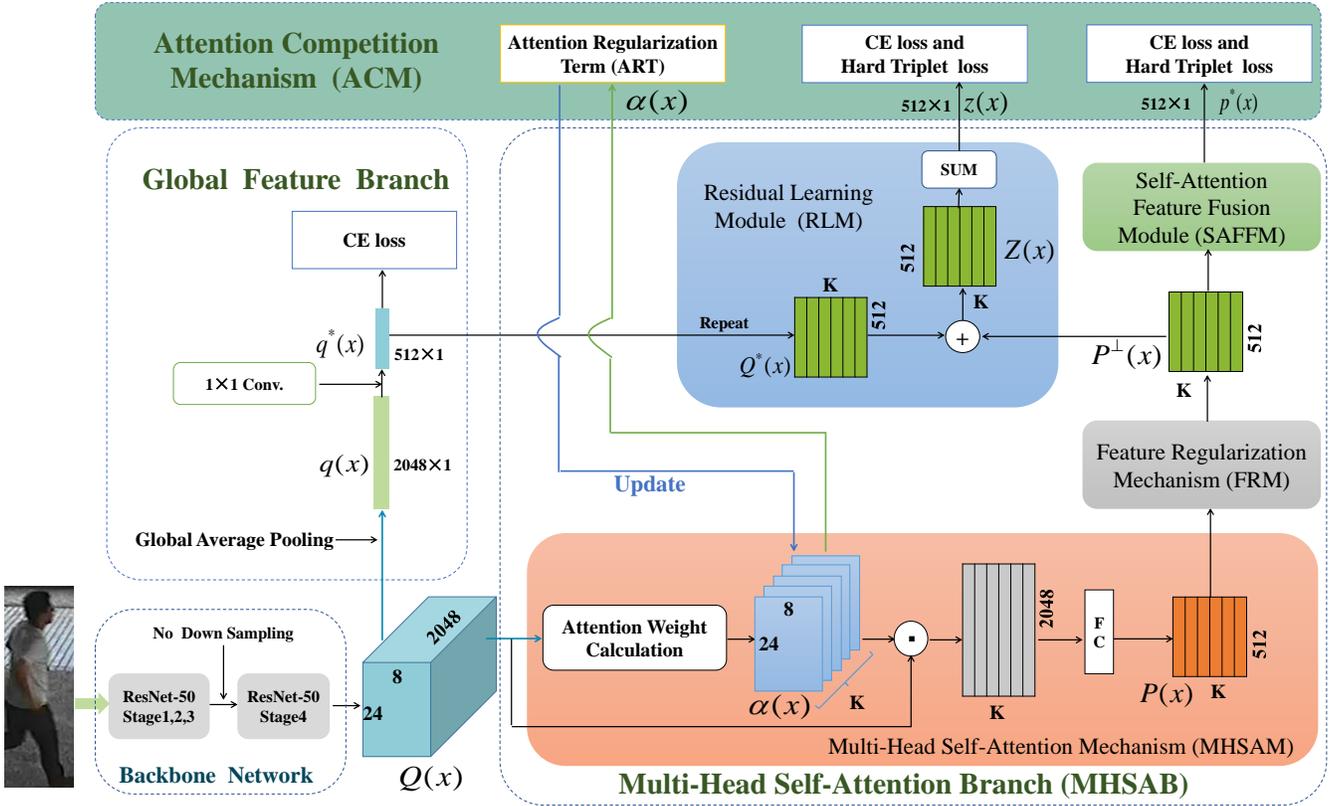


Fig. 2. The architecture of Multi-Head Self-Attention Network (MHSA-Net) for occlusion/standard person Re-ID task. The MHSA-Net contains three modules: the Global Feature Branch, the Multi-Head Self-Attention Branch (MHSAB), and the Attention Competition Mechanism (ACM). The CE loss denotes the cross entropy loss function.

additional model. Also, the MHSA-Net can more effectively capture unoccluded local information.

III. MHSA-NET OVERVIEW

Our MHSA-Net contains three modules: the Global Feature Branch, the Multi-Head Self-Attention Branch (MHSAB), and the Attention Competition Mechanism (ACM), as illustrated in Fig. 2.

The Global Feature Branch computes a basic large feature tensor $Q(x)$ and a global feature $q^*(x)$ for MHSAB and person matching (CE loss). We use the widely adopted Backbone Network ResNet-50 [21] to compute feature tensor $Q(x)$, then down-sample it to $q^*(x)$ for MHSAB and person matching.

The MHSAB, the core component in MHSA-Net, captures the key local information and outputs the fusion feature $p^*(x)$ for the person matching. The MHSAB contains four sub-modules: the (1) Multi-Head Self-Attention Mechanism (MHSAM), (2) Feature Regularization Mechanism (FRM), (3) Self-Attention Feature Fusion Module (SAFFM), and (4) Residual Learning Module (RLM). MHSAB outputs attention weights $\alpha(x)$, and **fusion features** $z(x)$ and $p^*(x)$ that capture key local information. These $\alpha(x)$, $z(x)$ and $p^*(x)$, will be refined in the Attention Competition Mechanism (ACM).

The ACM is composed of a series of loss functions and a regularization item; and it updates attention weights $\alpha(x)$, and **fusion features** $z(x)$ and $p^*(x)$ to enhance key person information and suppress non-key person information.

In the testing stage. For the standard person Re-ID task, we concatenate the feature vector $p^*(x) \in \mathbb{R}^{512}$ and $q^*(x) \in \mathbb{R}^{512}$ to find the best matching person in the gallery by comparing the squared distance, i.e. $d(a, b) = \|a - b\|_2^2$. For the occlusion person task, we only use the feature vector $p^*(x) \in \mathbb{R}^{512}$ to find the best matching person in the gallery by comparing the squared distance.

IV. GLOBAL FEATURE BRANCH (BASELINE)

Following recent state-of-the-art methods [77], [63], [5], [56], [47], [83], [19], we adopted ResNet-50 (pre-trained on ImageNet [26]) as the backbone network to encode a person image x . We modify the backbone ResNet-50 slightly to extract richer information via larger-sized high-level feature maps. The down-sampling operation at the beginning of stage 4 is not employed, then the output of the Backbone Network is $Q(x) \in \mathbb{R}^{24 \times 8 \times 2048}$. Following [83], [63], [19], we also append a series of downsampling operations to the large feature map $Q(x)$. As shown in Fig. 2, firstly we employ a global average pooling operation on the output feature $Q(x)$, the $24 \times 8 \times 2048$ tensor from the stage 4 of ResNet-50, to get a feature vector $q(x) \in \mathbb{R}^{2048}$. Then, $q(x)$ is further reduced to a 512-dimensional feature vector $q(x)^*$ through a 1×1 convolution layer, a batch normalization layer, and a **Rectified Linear Units (ReLU)** layer. Finally, the feature vector $q(x)^*$ is fed into the loss function, which is cross entropy loss \mathcal{L}_{CE}

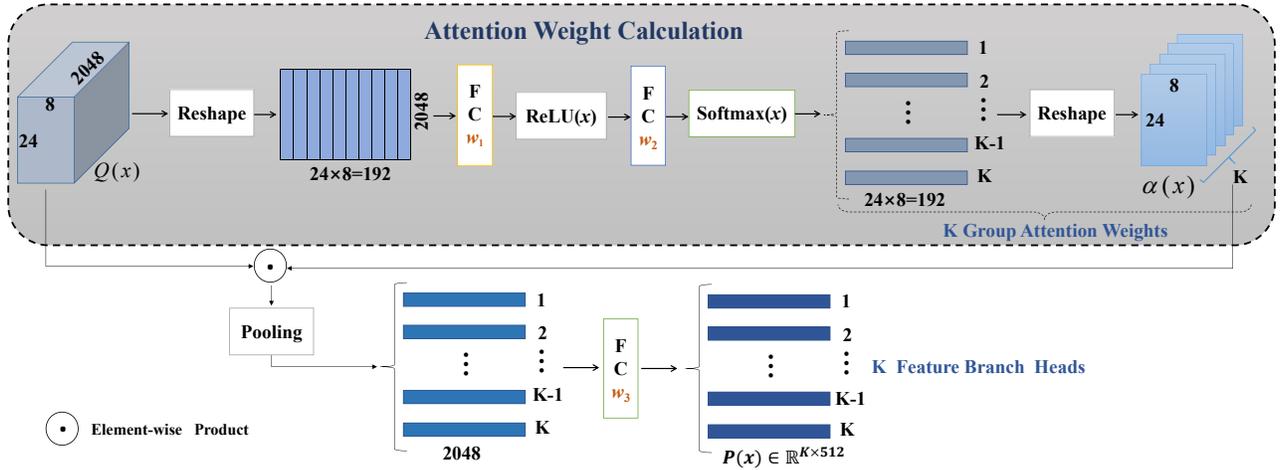


Fig. 3. The architecture of the Multi-Head Self-Attention Mechanism (MHSAM).

in this baseline model. The baseline model of our MHSA-Net is composed of the Global Feature Branch and the Backbone Network.

Unlike existing methods [77], [5], [56], [47], [73], we don't introduce the triplet loss into the GFB. In our experiments, we observed that incorporating triplet loss in the GFB negatively impacts the performance of MHSA-Net on generic person Re-ID with occlusions. It seems this more strict constraint on global features affects the local feature capturing in some degree. So, in our MHSA-Net, the loss function in the baseline model only contains \mathcal{L}_{CE} .

V. MULTI-HEAD SELF-ATTENTION BRANCH (MHSAB)

We introduce the Multi-Head Self-Attention Mechanism (MHSAM) [3], [25], [67] into the person Re-ID pipeline, to help the network capture key local information from occluded images. However, there are two issues need to be solved for this MHSAM in the occluded person Re-ID task.

(1) MHSAM [3], [25] can capture key local information using multiple embeddings; but these existing methods directly concatenate these embeddings, which result in a huge dimensional feature space, making search and training expensive and difficult. For our Re-ID task which is more complicated than the Natural Language Processing (NLP) task in [3], [78], we need a more effective design on MHSAM to output a low dimensional and efficient person descriptor for the person matching task.

(2) MHSAM produces multiple attention maps and feature embeddings for an image to encode rich information, which enhances the robustness of the deep model in representation learning [3]. But this design itself often makes different embeddings to redundantly encode similar or same personal information. Thus, it is desirable to make the generated embeddings diverse, namely, they capture various features of the person from different aspects.

Based on these observations, we propose a novel Multi-Head Self-Attention Branch (MHSAB) to tackle the above issues. MHSAB contains three components: the Multi-Head Self-Attention mechanism (MHSAM), Feature Regularization

Mechanism (FRM), and Self-Attention Feature Fusion Module (SAFFM). The MHSAM computes multiple attention maps for key sub-regions and multiple embeddings for each person image. The FRM contains a Feature Diversity Regularization Term (FDRT) and an Improved Hard Triplet Loss (IHTL) function. The FDRT enhances the diversity of the multiple embeddings in MHSAM, and the IHTL refines each individual embedding to better capture key information. The SAFFM adaptively combines multiple embeddings to produce a fused low-dimension feature vector.

A. Multi-Head Self-Attention Mechanism (MHSAM)

As described in Section I, it is desirable to adaptively capture key local features in unoccluded regions and avoid information from occluded regions. To achieve this, we adopt a Multi-Head Self-Attention mechanism (MHSAM) [3], [25], [34]. The architecture of MHSAM is shown in Fig. 3. Here, we build K -head for this MHSAM in two steps: (1) first, given a person image, we learn its K attention weights $\alpha(x) \in \mathbb{R}^{J \times K}$ (where $J = [24 \times 8]$) on each pixel $j \in J$ of feature maps $Q(x) \in \mathbb{R}^{J \times 2048}$. (2) Second, we compute the K attention-weighted embeddings of $Q(x)$ for this person image. Specifically:

(Step 1) Compute $\alpha(x)$ by the Attention Weight Calculation module (Figs. 3):

$$\alpha(x) = \text{softmax}(\omega_2 \text{ReLU}(\omega_1 Q(x)^T)), \quad (1)$$

where $Q(x) \in \mathbb{R}^{J \times 2048}$ is reshaped to a matrix in $\mathbb{R}^{192 \times 2048}$, $\alpha(x) \in \mathbb{R}^{K \times 192}$ is reshaped to a tensor in $\mathbb{R}^{24 \times 8 \times K}$, $\omega_2 \in \mathbb{R}^{K \times 512}$ and $\omega_1 \in \mathbb{R}^{512 \times 2048}$ are two parameter weight matrices to learn, and the softmax is applied pixel-wise so that on each pixel the K attention weights sum up to one.

(Step 2) Multiply the attention weight $\alpha(x)$ with feature maps $Q(x)$, and further apply a non-linear transformation, to get K attention-weighted embeddings $P \in \mathbb{R}^{K \times 512}$ (Figs. 3):

$$P(x) = \text{AvgPool}(\alpha(x) \odot Q(x))\omega_3 + b_3, \quad (2)$$

where $\omega_3 \in \mathbb{R}^{2048 \times 512}$ is the parameter weight matrix, $\text{AvgPool}(\cdot)$ is the average pooling operation, and $b_3 \in \mathbb{R}^{512}$ is

the bias to learn for the fully connection layer “FC” in Fig. 3. Since this, we can obtain K feature branch heads, and the number of the heads is K . The \odot is the element-wise product operation.

The attention weights $\alpha(x)$ in Eq. (1) are adaptively learned toward the objective of person matching in Re-ID. Greater α values indicate bigger importance of pixels/local regions and vice versa. As some examples shown in Figs. 1 and 7, key information from unoccluded regions can be captured by MHSAM, while occluded regions can be suppressed. Here, the hyper-parameter K is discussed in the Subsection VII-E1.

B. Feature Regularization Mechanism (FRM)

FRM contains a Feature Diversity Regularization Term (FDRT) and an Improved Hard Triplet Loss (IHTL). The FDRT encourages the multiple embeddings $P(x)$ to cover more key local information from various respects. The IHTL refines the embeddings so that they individually can better serve person matching. FRM takes in $P(x) \in \mathbb{R}^{K \times 512}$, and outputs a new tensor $P^\perp(x) \in \mathbb{R}^{K \times 512}$.

1) Feature Diversity Regularization Term (FDRT):

The K embeddings directly produced by MHSAM tend to capture similar/same person information redundantly. To avoid this, following [3], we also introduce the Feature Diversity Regularization Term (FDRT) into MHSAM, to regularize the K representations and enforce their diversity.

The K embeddings in MHSAM are not overcomplete [23], [13]. So we can restrict the Gram matrix of K embeddings to be close to an identity matrix under Frobenius norm. Firstly, we create a Gram matrix $G(x)$ of $P(x)$ by $G(x) = P(x)P(x)^T$. Each element in $G(x)$ denotes the correlation between $P(x)$. Here, $P(x)$ is normalized so that they are on an L_2 ball. Secondly, to enhance the diversity of the K feature vectors in $P(x)$, we minimize the deviation of $G(x)$ from the identity matrix. Therefore, we define the Feature Diversity Regularization Term (FDRT) as

$$\mathcal{L}_{FDRT} = \frac{1}{K^2} \|G(x) - I\|_1, \quad (3)$$

where $G(x)$ is the gram matrices of $P(x)$, and $I \in \mathbb{R}^{K \times K}$ is an identity matrix. With FDRT, the K embeddings $P(x)$ are more diverse and can capture key information from different perspectives, which enhances the model robustness.

2) *Improved Hard Triplet Loss (IHTL)*: MHSAM produces K embeddings $P(x) \in \mathbb{R}^{K \times H}$ for each person image. To further filter out non-key information, we design a new loss function to help train the network so that each individual embedding can be used separately for person matching. We are inspired by the hard triplet loss [1], which uses a hard sample mining strategy to achieve desirable performance. Hence, we propose an Improved Hard Triplet Loss (IHTL) by revising the hard triplet loss [1].

Before defining the Improved Hard Triplet Loss (IHTL), we firstly organize the training samples into a set of triplet feature units, $S = (s(x^a), s(x^p), s(x^n))$, or simply $S = (s^a, s^p, s^n)$ in the following. The raw person image triplet units is $X = (x^a, x^p, x^n)$. Here, (s^a, s^p) represents a positive

pair of features $y^a = y^p$, and (s^a, s^n) indicates a negative pair of features with $y^a \neq y^n$. Here, $y \in Y$ is the person ID.

In the hard triplet loss [1], a hard-sample mining strategy is introduced: a positive sample pair with the largest distance is defined as the *hard positive sample pair*; the negative sample pair with the smallest distance is defined as the *hard negative sample pair*. The hard triplet loss function can then be defined using hard sample pairs:

$$\mathcal{T}_{HardTriplet} = \ln(1 + \exp(\max_{x^a, x^p} d(s^a, s^p) - \min_{x^a, x^n} d(s^a, s^n))), \quad (4)$$

Based on the hard triplet loss function, we define an Improved Hard Triplet Loss (IHTL). We define the *improved hard positive sample pair* and *improved hard negative sample pair* in two steps: **(I)**: Between each sample image pair, $K \times K$ distances can be computed, because each person image has K embeddings $P(x) \in \mathbb{R}^{K \times 512}$ in MHSAM. We use the largest distance from these distances to measure the embeddings of the positive sample pairs, and use the smallest distance from these distances for the negative sample pairs.

(II): We further use the hard samples mining strategy [1] to define the hard sample pairs. The improved hard positive sample pair is $\max_{x^a, x^p} \max_{i, j} d(P(x^a)_i, P(x^p)_j)$; the improved hard negative sample pair is $\min_{x^a, x^n} \min_{i, j} d(P(x^a)_i, P(x^n)_j)$. The Improved Hard Triplet loss is defined as:

$$\mathcal{T}_{IHTL} = \ln(1 + \exp(\max_{x^a, x^p} \max_{i, j} d(P(x^a)_i, P(x^p)_j) - \min_{x^a, x^n} \min_{i, j} d(P(x^a)_i, P(x^n)_j))), \quad (5)$$

where $i, j \in \{1, 2, \dots, K\}$, $d(a, b) = \|a - b\|_2^2$ denotes the squared distance in feature space. Here, During training, the IHTL refines embeddings so that they individually can perform better person matching. This encourages the embeddings to focus on important information.

C. Self-Attention Feature Fusion Module (SAFFM)

The output of **FDRT**, the K embeddings $P^\perp(x) \in \mathbb{R}^{K \times 512}$ covers various properties of a person image. But directly using $P^\perp(x)$ by concatenation will lead to dimension explosion in person matching. Thus, we design a Self-Attention Feature Fusion Module (SAFFM) to first learn K attentional weights by a series of neural networks, then fuse $P^\perp(x)$ to get a lower-dimensional $p(x)^* \in \mathbb{R}^{512}$.

Specifically, **Step-1**, compute the attentional weight $\beta(x) \in \mathbb{R}^{K \times 512}$ (Fig. 2 and Fig. 4). The matrix $P^\perp(x) \in \mathbb{R}^{K \times 512}$ is transposed to $P^*(x) \in \mathbb{R}^{512 \times K}$, then compute β by

$$\beta(x) = \text{softmax}(\omega_5 \text{ReLU}(\omega_4 P^*(x))), \quad (6)$$

where $\omega_4 \in \mathbb{R}^{512 \times 1024}$ and $\omega_5 \in \mathbb{R}^{1024 \times 512}$ are two parameter weight matrices to learn, and the softmax is applied pixel-wise so that each pixel on the each attention vector of the $\beta(x)$ sum up to one.

Step-2, compute the self-attention weighted feature vector $p(x)^* \in \mathbb{R}^{512}$, by

$$p(x)^* = \sum_{i=1}^K [\beta(x) \odot P^\perp]_i. \quad (7)$$

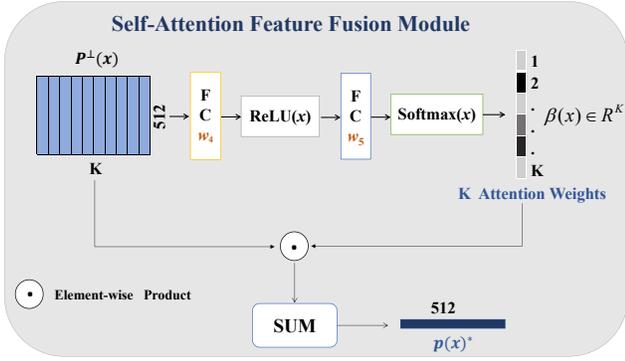


Fig. 4. The Self-Attention Feature Fusion Module.

Here, $p(x)^* \in \mathbb{R}^{512}$, \odot is the element-wise product operation.

SAFFM reduces the dimension of the multiple embeddings $P^\perp(x)$ for both training and testing. In the training stage, $p^*(x)$ is also fed to the cross entropy loss and the hard triplet loss function,

$$\mathcal{L}_{SAFFM} = \mathcal{L}_{CE}^* + \mathcal{T}_{HardTriplet}^* \quad (8)$$

Here, the input to both \mathcal{L}_{CE}^* and $\mathcal{L}_{HardTriplet}^*$ is $p^*(x)$.

D. Residual Learning Module

As shown in the RLM module in Fig. 2, with MHSAM, $P^\perp(x)$ aims to capture key information in unoccluded regions from local perspective; while $q(x)^*$ captures global information of the whole person image. To prevent $P^\perp(x)$ from being redundant with $q(x)^*$, we cast their feature fusion as a residual learning task. Specifically, (1) to match with the dimension $K \times 512$ of $P^\perp(x)$, we copy $q(x)^*$ for K times to obtain $Q(x)^* \in \mathbb{R}^{K \times 512}$. (2) The input to the residual block includes global feature $Q(x)^*$ and local feature $P^\perp(x)$. The parameters $(\omega_1, \omega_2, \omega_3, b_3)$ of $P^\perp(x)$ will be optimized. (3) We define the residual learning embedding as

$$Z(x) = Norm(Q(x)^* + P^\perp(x)), \quad (9)$$

where $Norm(\cdot)$ denotes the layer normalization [31]. This RLM encourages $P^\perp(x)$ to only capture important local information.

In the training stage, $Z(x) \in \mathbb{R}^{K \times 512}$ is simply summed along the first dimension to obtain $z(x) \in \mathbb{R}^{512}$. And $z(x)$ is also fed into the cross entropy loss and the hard triplet loss function, i.e.

$$\mathcal{L}_{ReN} = \mathcal{L}_{CE}^{**} + \mathcal{L}_{HardTriplet}^{**} \quad (10)$$

Here, the input of \mathcal{L}_{CE}^{**} and $\mathcal{L}_{HardTriplet}^{**}$ is $z(x)$. And $z(x)$ does not participate in person matching in the testing stage.

Finally, the **loss functions in MHSAB** are summarized as

$$\mathcal{L}_{MHSAB} = \mathcal{L}_{SAFFM} + \lambda_1 \mathcal{L}_{FDRT} + \mathcal{L}_{ReN} + \lambda_2 \mathcal{T}_{IHTL}, \quad (11)$$

where λ_1 and λ_2 are the balance parameters (see detail in Subsection VII-E2 and the Subsection VII-E3).

VI. ATTENTION COMPETITION MECHANISM

MHSAB enhances attention on key sub-regions, but the extracted attention maps still contain some non-key information. We propose an Attention Competition Mechanism (ACM) to further refine the attention weights.

In [22], an attention competition strategy was proposed to filter out attention information of the non-key words in the Text-to-Image generation task. This idea was composed of an attention regularization term and a series of cross-modal matching loss functions. This has been shown effective in the Text-to-Image generation task. In the image generation: an attention regularization term can effectively filter out the attention information of non-key words; the cross-modal matching loss functions can effectively enhance or preserve the attention information of the key words according to the objective. Similarly, we believe it can also help the person Re-ID model filter out the attention information of the non-key sub-regions from the person images. Therefore, we also design a similar strategy in this Person Re-ID pipeline. To our knowledge, this is the first time a competition strategy was designed for Re-ID task. Through a series of experiments, we observe that this mechanism is promising.

Specifically, we use an attention regularization term to suppress non-key information, and use the aforementioned person Re-ID loss function \mathcal{L}_{MHSAB} to enhance attention on important regions. The attention regularization term [22] is defined as:

$$\mathcal{L}_C = \sum_{i,j} (\min(\alpha_{i,j}, \gamma))^2, \quad (12)$$

where the subscript ‘‘C’’ stands for ‘‘Competition’’, and $\gamma > 0$ is a threshold. Fig. 5 shows a schematic diagram of the ACM. The grey columns illustrate attention weights on non-key sub-regions, and the green columns are for weights on key regions. In the initial state of training, as shown in subfig (a), all attention weights in α are small. In ACM, the attention regularization term \mathcal{L}_C sets a threshold and pushes the attention weights lower than this threshold toward zero; while \mathcal{L}_{MHSAB} increases attention weights of sub-regions if they benefit person matching. An illustration of this procedure is given in (b).

The total loss functions in MHSAB-Net. The total loss function \mathcal{L}_{Total} is

$$\mathcal{L}_{Total} = \mathcal{L}_{MHSAB} + \mathcal{L}_{GFB} + \lambda_3 \mathcal{L}_C, \quad (13)$$

where λ_3 is the balance parameter (see its discussion in the Section VII-F).

VII. EXPERIMENT

To evaluate the MHSAB-Net, we conduct extensive experiments on three widely used generic person Re-ID benchmarks, i.e. **Market-1501** [38], **DukeMTMC-reID** [46], [74] and **CUHK03** [53] datasets, and four occluded person Re-ID benchmarks, i.e. **Occluded-DukeMTMC** [29], **P-DukeMTMC-reID** [30], **Partial-REID** [55] and **Partial-iLIDS** [39]. First, we compare the performance of MHSAB-Net with state-of-the-art methods on these datasets. Second,

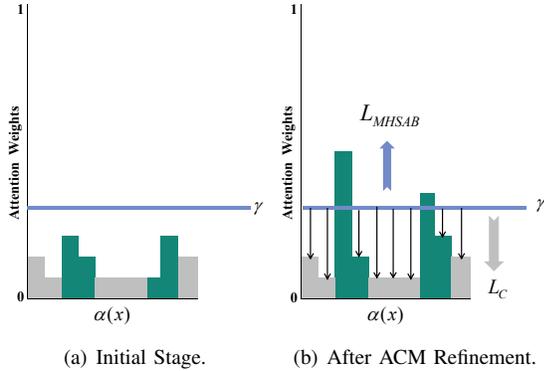


Fig. 5. The schematic diagram of the attention competition process on $\alpha(x)$. Grey columns are the attention weights of non-key sub-regions, and green columns are the weights of key sub-regions.

we perform ablation studies to validate the effectiveness of each component.

A. Datasets and Evaluation

We follow almost all person Re-ID approaches [81], [77], [83], [63], [46], [74], [39], [55], [29], [30], [4], [70] to set the following seven person Re-ID datasets.

Market-1501 [38] has 32,668 labeled images of 1,501 identities collected from 6 camera views. The dataset is partitioned into two non-overlapping parts: the training set with 12,936 images from 751 identities, and the test set with 19,732 images from 750 identities. In the testing stage, 3,368 query images from 750 identities are used to retrieve the persons from the rest of the test set, i.e. the gallery set.

DukeMTMC-reID [46], [74] is another large-scale person Re-ID dataset. It has 36,411 labeled images of 1,404 identities collected from 8 camera views. The training set consists of 16,522 images from 702 identities; We used 2,228 query images from the other 702 identities, and 17,661 gallery images.

CUHK03 [53] is a challenging Re-ID benchmark. It has 14,096 images of 1,4674 identities captured from 6 cameras. It contains two datasets. **CUHK03-Labeled**: the bounding boxes of person images are from manual labeling. **CUHK03-Detected**: the bounding boxes of person images are detected from deformable part models (DPMs), which is more challenging due to severe bounding box misalignment and background cluttering. Following [81], [77], [83], [63], we used the 767/700 split [53] of the detected images.

Occluded-DukeMTMC [29] has 15,618 training images, 17,661 gallery images, and 2,210 occluded query images. We use this dataset to evaluate our MHSA-Net in Occluded Person Re-ID task.

P-DukeMTMC-reID [30] is a modified version based on DukeMTMC-reID [46], [74]. There are 2,652 images (665 identities) in the training set, 2,163 images (634 identities) in the query set and 9,053 images in the gallery set.

Partial-REID [55] is a specially designed partial person Re-ID benchmark that has 600 images from 60 people. Each person has five partial images in query set and five full-body images in gallery set. These images are collected at

a university campus under different viewpoints, backgrounds, and occlusions.

Partial-iLIDS [39] is a simulated partial person Re-ID dataset based on the iLIDS dataset. It has a total of 476 images of 119 people.

Evaluation Protocol. We employed two standard metrics adopted in most person Re-ID approaches, namely, the cumulative matching curve (CMC) that generates ranking accuracy, and the mean Average Precision (mAP). The CMC curve shows the probability that a query identity appears in different-sized candidate lists. This evaluation measurement is valid only if there is only one ground truth match for a given query. In this paper, we report the Rank-1 accuracy. The mAP calculates the area under the Precision-Recall curve, which is known as average precision (AP). Then, the mean value of APs of all queries, i.e., mAP, is calculated, which considers both precision and recall of an algorithm, thus providing a more comprehensive evaluation.

B. Implementation Details

Following many recent approaches [29], [77], [63], [83], [9], the input images are re-sized to 384×128 and then augmented by random horizontal flip and normalization in the training stage. In the testing stage, the images are also re-sized to 384×128 and augmented only by normalization. Using the ImageNet pre-trained ResNet-50 as the backbone, our network is end-to-end in the whole training stage. Our network is trained using 2 GTX 2080Ti GPUs with a batch size of 128. Each batch contains 32 identities, with 4 samples per identity. We use Adam optimizer [36] with 400 epochs. The base learning rate is initialized to 10^{-3} with a linear warm-up [44] in first 50 epochs, then decayed to 10^{-4} after 200 epochs, and further decayed to 10^{-5} after 300 epochs. The whole training procedure has 400 epochs and takes approximately 2 hours. Our MHSA-Net achieves the satisfactory performance in the general person Re-ID and occluded person Re-ID tasks, when $\lambda_1 = 1e-4$, $\lambda_2 = 1.0$, $\lambda_3 = 1e-3$, $\gamma = 1e-3$ and $K = 8$.

C. Comparison with state-of-the-art Methods

In this subsection, we compared MHSA-Net with a series of state-of-the-art approaches on seven person Re-ID datasets. Here, MHSA-Net concatenates the local feature $p(x)^*$ and global feature $q(x)^*$ to conduct the person matching task. Compared with the proposed MHSA-Net, MHSA-Net[†] indicates that we drop the \mathcal{L}_{CE} in the training process. The MHSA-Net* only uses the local feature $p(x)^*$ to conduct the person Re-ID task.

Person Re-ID on General Datasets. Firstly, we compared MHSA-Net with the state-of-the-art generic person Re-ID approaches on Market-1501, DukeMTMC-Re-ID, CUHK03-Labeled, and CUHK03-Detected datasets, and reported the results in Tables I. We randomly set $K = 5$ and $K = 8$ in these experiments (Through experiments, we observed that the K value does not affect the result much. Some discussions on different K values are given in Subsection VII-E1). Our MHSA-Net gets *Rank-1* = 94.6, 87.3, 73.4, 75.8 and *mAP* = 84.0, 73.1, 70.2, 73.0 for Market-1501, DukeMTMC-reID,

TABLE I

THE COMPARISON WITH THE MANY STATE-OF-THE-ART GENERIC PERSON RE-ID METHODS ON MARKET-1501, DUKEMTMC-reID, AND CUHK03 DATASETS. MHSA-NET[†] INDICATES THAT WE DROP THE \mathcal{L}_{CE} IN THE TRAINING PROCESS. MHSA-NET* INDICATES THAT WE ONLY USE THE KEY LOCAL FEATURE $p^*(x)$ FROM MHSAB+ACM TO CONDUCT THE PERSON MATCHING TASK.

Method	Market-1501		DukeMTMC-reID		CUHK03-Detected		CUHK03-Labeled	
	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP
MLFN [59]	90.0	74.3	81.0	62.8	52.8	47.8	54.7	49.2
HA-CNN [54]	91.2	75.7	80.5	63.8	41.7	38.6	44.4	41.0
PCB+RPP[63]	93.8	81.6	83.3	69.2	62.8	56.7	-	-
Mancs[7]	93.1	82.3	84.9	71.8	65.5	60.5	-	-
PAN [75]	82.8	63.4	71.6	51.5	36.3	34.0	36.9	35.0
FANN[48]	90.3	76.1	-	-	69.3	67.2	-	-
VCFL[16]	90.9	86.7	-	-	70.4	70.4	-	-
PGFA [29]	91.2	76.8	82.6	65.5	-	-	-	-
HACNN+DHA-NET [52]	91.3	76.0	81.3	64.1	-	-	-	-
IANet[47]	94.4	83.1	87.1	73.4	-	-	-	-
BDB[83]	94.2	84.3	86.8	72.1	72.8	69.3	73.6	71.7
AANet[9]	93.9	83.4	87.7	74.3	-	-	-	-
CAMA[56]	94.7	84.5	85.8	72.9	66.6	64.2	-	-
OSNet[77]	94.8	84.9	88.6	73.5	72.3	67.8	-	-
RANgEv2[57]	94.7	86.8	87.0	78.2	64.6	61.6	67.4	64.3
JWSAA[43]	94.8	83.2	88.3	75.6	72.3	67.8	-	-
HOReID [18]	94.2	84.9	86.9	75.6	-	-	-	-
SCSN (4-stage) [6]	92.4	88.3	91.0	79.0	84.7	81.0	86.8	84.0
SCSN (3-stage) [6]	95.7	88.5	90.1	79.0	84.1	80.2	86.3	83.3
RGA-SC [66]	95.8	88.1	86.1	74.9	77.3	73.3	80.4	76.4
Baseline	92.0	78.8	81.0	62.8	56.3	53.0	58.6	55.2
MHSA-Net (K=5)	94.3	83.5	87.1	73.0	73.4	70.2	75.8	73.0
MHSA-Net (K=8)	94.6	84.0	87.3	73.1	72.8	69.3	75.6	72.7
MHSA-Net* (K=8)	94.0	82.9	86.3	72.5	72.4	69.7	74.4	72.0
MHSA-Net [†] (K=8)	94.3	82.5	87.0	72.6	72.7	69.9	75.2	72.3
MHSA-Net+Re-ranking [79] (K=8)	95.5	93.0	90.7	87.2	80.2	80.9	82.6	84.2

TABLE II

THE COMPARISON WITH THE OTHER OCCLUDED PERSON RE-ID METHODS ON OCCLUDED-DUKEMTMC DATASET. MHSA-NET[†] INDICATES THAT WE DROP THE \mathcal{L}_{CE} IN THE TRAINING PROCESS. MHSA-NET* INDICATES THAT WE ONLY USE THE KEY LOCAL FEATURE $p^*(x)$ FROM MHSAB+ACM TO CONDUCT THE PERSON MATCHING TASK. THE **FIRST**, **SECOND** AND **THIRD** HIGHEST SCORES ARE SHOWN IN **RED**, **GREEN** AND **BLUE** RESPECTIVELY.

Method	Occluded-DukeMTMC			
	Rank-1	Rank-5	Rank-10	mAP
Random Erasing [80]	40.5	59.6	66.8	30.0
HA-CNN [54]	34.4	51.9	59.4	26.0
Adver Occluded [24]	44.5	-	-	32.2
PCB [63]	42.6	57.1	62.9	33.7
Part Bilinear [65]	36.9	-	-	-
FD-GAN [64]	40.8	-	-	-
DSR [39]	40.8	58.2	65.2	30.4
SFR [41]	42.3	60.3	67.3	32.0
PGFA [29]	51.4	68.6	74.9	37.3
HOReID [18]	55.1	-	-	43.8
PVPM+Aug [17]	57.3	72.6	77.2	45.7
Baseline	38.9	53.5	60.1	25.6
MHSA-Net* (K=8)	59.7	74.3	79.5	44.8
MHSA-Net (K=8)	55.4	70.2	76.4	42.4
MHSA-Net [†] (K=8)	58.2	73.2	78.4	43.1

CUHK03-Detected and CUHK03-Labeled, respectively. If we introduce the Re-ranking [79] into the MHSA-Net, i.e. MHSA-Net+Re-ranking (K=8), the accuracy further increases to *Rank-1*= 95.5 , 90.7 , 80.2 , 82.6 and *mAP*= 93.0 , 87.2 , 80.9 , 84.2 for *Market-1501*, *DukeMTMC-reID*, *CUHK03-Detected* and *CUHK03-Labeled*, respectively. Recently, the state-of-the-art performance on *Market-1501* and *DukeMTMC-Re-ID* has been saturated. Yet the MHSA-Net still gains effective

TABLE III

THE COMPARISON WITH THE OTHER OCCLUDED PERSON RE-ID METHODS ON P-DUKEMTMC-reID DATASET. MHSA-NET[†] INDICATES THAT WE DROP THE \mathcal{L}_{CE} IN THE TRAINING PROCESS. MHSA-NET* INDICATES THAT WE ONLY USE THE KEY LOCAL FEATURE $p^*(x)$ FROM MHSAB+ACM TO CONDUCT THE PERSON MATCHING TASK. THE **FIRST**, **SECOND** AND **THIRD** HIGHEST SCORES ARE SHOWN IN **RED**, **GREEN** AND **BLUE** RESPECTIVELY.

Method	P-DukeMTMC-reID			
	Rank-1	Rank-5	Rank-10	mAP
OSNet[77]	33.7	46.5	54.0	20.1
PCB+RPP[63]	40.4	54.6	61.1	23.4
PCB[63]	43.6	57.1	63.3	24.7
PGFA [29]	44.2	56.7	63.0	23.1
PVPM+Aug [17]	51.5	64.4	69.6	29.2
Baseline	61.0	72.5	78.4	27.0
MHSA-Net* (K=8)	70.7	81.0	84.6	41.1
MHSA-Net (K=8)	67.9	79.7	83.7	37.6
MHSA-Net [†] (K=8)	69.6	81.4	85.0	37.5

improvement over the baseline model and outperforms most existing methods. The CUHK03 is the most challenging dataset among the three. Following the data setting in [81], [77], [83], [63], MHSA-Net also outperforms the most state-of-the-art methods on both CUHK03-Labeled and CUHK03-Detected datasets.

Occluded Person Re-ID. A feature of MHSA-Net is that it handles Re-ID of occluded persons well. So, we also compared MHSA-Net with a series of occluded person re-id methods on the Occluded-DukeMTMC dataset, P-DukeMTMC-reID dataset, Partial-REID dataset, and Partial-iLIDS dataset.

Occluded-DukeMTMC and P-DukeMTMC-reID. MHSA-Net* only uses the local features $p(x)^*$ for the

TABLE IV

THE COMPARISON WITH THE OTHER OCCLUDED PERSON RE-ID METHODS ON PARTIAL-REID AND PARTIAL-IIDS. MHSA-Net[†] INDICATES THAT WE DROP THE \mathcal{L}_{CE} IN THE TRAINING PROCESS. MHSA-Net* INDICATES THAT WE ONLY USE THE KEY LOCAL FEATURE $p^*(x)$ FROM MHSAB+ACM TO CONDUCT THE PERSON MATCHING TASK. THE FIRST, SECOND AND THIRD HIGHEST SCORES ARE SHOWN IN RED, GREEN AND BLUE RESPECTIVELY.

Method	Partial-REID		Partial iLIDS	
	Rank-1	Rank-3	Rank-1	Rank-3
MTRC [49]	23.7	27.3	17.7	26.1
AMC+SWM [55]	37.3	46.0	21.0	32.8
DSR [39]	50.7	70.0	58.8	67.2
SFR [41]	56.9	78.5	63.9	74.8
FPR [40]	81.0	-	68.1	-
PGFA [29]	68.0	80.0	69.1	80.9
PVPM+Aug [17]	80.6	84.2	68.7	81.4
HOReID [18]	85.3	91.0	72.6	86.4
Baseline	68.8	81.7	66.4	79.0
MHSA-Net* (K=8)	85.7	91.3	74.9	87.2
MHSA-Net [†] (K=8)	85.5	91.0	74.1	86.6
MHSA-Net (K=8)	81.3	87.7	73.6	85.4

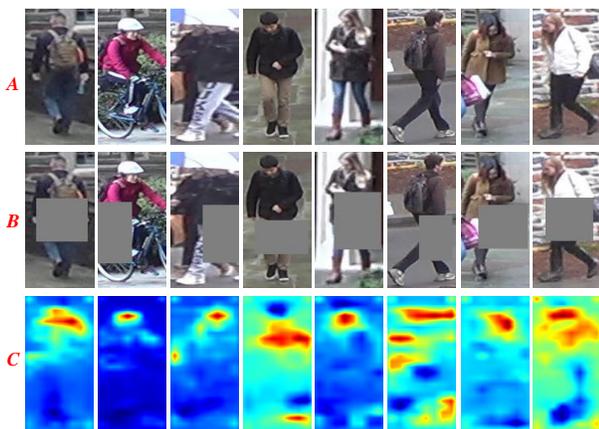


Fig. 6. Visualization of feature map corresponding to random manual occlusion. Image Group A (the first row) is the normal person images. We randomly erase part of Image Group A to get Image Group B (the second row). The Image Group C (the third row) is corresponding feature maps of Image Group B.

occluded Re-ID task. As shown in Table II, on Occluded-DukeMTMC, our MHSA-Net* achieves **59.7 Rank-1 accuracy and 44.8 mAP**, which outperforms most previous methods. Compared with the baseline model, the MHSA-Net* gains 20.8 Rank-1 and 19.6 mAP improvement. As shown in Table III, on P-DukeMTMC-reID, our MHSA-Net* achieves **70.7 Rank-1 accuracy and 41.1 mAP**, which outperforms all the previous methods. Compared with the baseline model, the MHSA-Net* gains 9.7 Rank-1 and 14.1 mAP improvement.

The MHSA-Net combines both global feature $q(x)^*$ and local feature $p(x)^*$ to do the occluded person Re-ID task. The global feature $q(x)^*$ captures global information from the whole person image. Hence, it inevitably encodes contents of scene regions that occlude the person, and this leads to decreased performance. This can be remedied by reducing the constraints on the global feature branch. Specifically, if we drop the \mathcal{L}_{CE} in the Global Feature Branch (GFB), the extraction of global feature becomes a simple downsampling

from local features, making this global feature $q(x)^*$ less sensitive to occlusions. We denote the pipeline using such a design as MHSA-Net[†]. In Table II and Table III, the performance of MHSA-Net[†] is clearly better than MHSA-Net, and only slightly worse than MHSA-Net*. The MHSA-Net[†] also achieves a competitive performance in Table I.

In summary, for general database, we can use MHSA-Net. For datasets with certain occlusions, we can use MHSA-Net[†], which best balances the global and local features. For datasets with severe occlusions, we can use MHSA-Net*, where local features play more important roles. As shown in Fig. 6, for the manually drawn occlusion, the MHSAM can better avoid the feature extraction of the occlusion part, and better extract the key person information of the non-occlusion part.

Partial-REID and Partial-iLIDS. The comparison of Re-ID on these two datasets is shown in Table IV. We also trained our model using the Market-1501 training set. Our MHSA-Net* and MHSA-Net[†] also achieve the best performance on both datasets. In both of these two data settings, compared with the baseline model, MHSA-Net* and MHSA-Net[†] gain a large improvement on both datasets. Like in Occluded-DukeMTMC and P-DukeMTMC-reID, MHSA-Net* and MHSA-Net[†] have better performance than MHSA-Net in Partial-REID and Partial-iLIDS datasets.

D. Ablation Study of MHSA-Net

We conducted ablation studies to show effectiveness of each component in the MHSA-Net. We show the ablation experiments results in Tables V and VI. By a series of discussions of hyper-parameters in Subsection VII-E, we found the best hyper-parameters setting for the proposed MHSA-Net: $\lambda_1 = 1e - 4$ and $\lambda_2 = 1.0$ in Eq. 11; $\lambda_3 = 1e - 3$ in Eq. 13; $\gamma = 1e - 3$ in Eq. 12. In these experiments, we set $K = 8$ in MHSAM, as we observed it produces stable and effective person matching results.

Through ablation studies, we have the following observations: (1) Each individual component effectively improves the performance of the baseline model, as shown in Table V. Compared with the baseline model, the entire MHSA-Net (K=8) achieves: 2.6 Rank-1 and 7.2 mAP improvement on Market-1501; 6.3 Rank-1 and 10.3 mAP improvement on DukeMTMC-Re-ID; 16.5 Rank-1 and 16.3 mAP on CUHK03-Detected; 17.0 Rank-1 and 17.5 mAP on CUHK03-Labeled. When $K = 5$ or 6, the change in performance is minor. (2) We show the effectiveness of each embedding in $P^\perp(x) \in \mathbb{R}^{K \times 512}$, i.e. $P^\perp(x)_i, i = 1, 2, \dots, K, K = 8$. We only use the $P^\perp(x)_i$ to search the target person. Here, we conduct the ablation studies on Occluded-DukeMTMC and CUHK03-Detected datasets in Table VI. As we can see the Table VI, compared with the baseline model, each feature $P^\perp(x)_i$ achieves large improvement over the baseline model in these two datasets. So, it indicates that each feature in $P^\perp(x)$ can better carry on both the *generic* and *occluded* person Re-ID tasks.

In the visualization results in Fig. 7, compared with Baseline, MHSAB can effectively capture more key sub-regions.

TABLE V

RESULTS PRODUCED BY COMBINING DIFFERENT COMPONENTS OF THE MHSA-NET. MHSA-NET[†] INDICATES THAT WE DROP THE \mathcal{L}_{CE} IN THE TRAINING PROCESS. MHSA-NET* = MHSA+IHTL+FDRT+ACM DENOTES THAT WE ONLY USE THE LOCAL FEATURE $p^*(x)$ CONDUCTS PERSON RE-ID TASK. FROM “BASELINE” TO “MHSA-NET[†]” IN THIS TABLE, THE PARAMETER K IS SET TO 8 IN IMPLEMENTATION.

Method	Market-1501		DukeMTMC-reID		CUHK03-Detected		CUHK03-Labeled	
	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP
Baseline	92.0	78.8	81.0	62.8	56.3	53.0	58.6	55.2
Baseline+MHSAM (K=8)	93.0	80.2	85.2	70.2	68.1	64.4	72.1	69.4
Baseline+MHSAM+ACM (K=8)	93.4	81.9	86.5	73.1	69.9	65.8	72.9	70.1
Baseline+MHSAM+FDRT (K=8)	93.6	82.1	86.3	72.6	70.0	65.7	73.4	69.5
Baseline+MHSAM+IHTL (K=8)	93.5	82.2	86.4	72.6	71.2	67.6	73.6	70.5
Baseline+MHSAM+IHTL+ACM (K=8)	94.1	83.3	86.8	72.7	72.2	69.9	75.9	73.4
Baseline+MHSAM+FDRT+ACM (K=8)	94.1	83.6	86.5	73.1	71.1	69.4	72.9	68.9
Baseline+MHSAM+IHTL+FDRT (K=8)	93.9	83.2	86.9	73.2	72.2	69.3	75.0	72.3
MHSA-Net* (K=8)	94.0	82.9	86.3	72.5	72.4	69.7	74.4	72.0
MHSA-Net [†] (K=8)	94.3	82.5	87.0	72.6	72.7	69.9	75.2	72.3
MHSA-Net (K=5)	94.3	83.5	87.1	73.0	73.4	70.2	75.8	73.0
MHSA-Net (K=6)	94.2	84.1	87.0	73.0	73.4	70.1	75.2	72.8
MHSA-Net (K=8)	94.6	84.0	87.3	73.1	72.8	69.3	75.6	72.7

TABLE VI

RESULTS ON OCCLUDED-DUKEMTMC AND CUHK03-DETECTED DATASET FOR EACH EMBEDDING IN THE $P^\perp(x)$. THE BOLD IS THE BEST RESULT.

Method	Occluded-DukeMTMC		CUHK03-Detected	
	Rank-1	mAP	Rank-1	mAP
Baseline	38.9	25.6	56.3	53.0
$P^\perp(x)_1$	55.2	40.4	68.7	65.0
$P^\perp(x)_2$	55.7	40.4	68.1	64.4
$P^\perp(x)_3$	53.3	39.6	69.4	66.4
$P^\perp(x)_4$	54.9	40.1	68.6	64.7
$P^\perp(x)_5$	54.5	39.9	69.1	66.3
$P^\perp(x)_6$	54.3	40.2	69.4	66.4
$P^\perp(x)_7$	53.2	38.7	70.0	67.7
$P^\perp(x)_8$	53.8	39.4	68.9	65.7

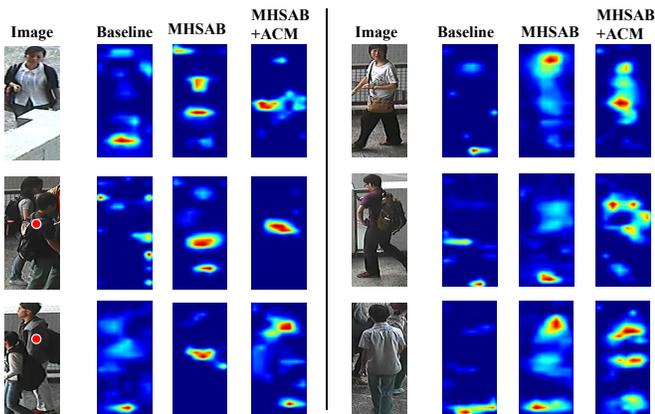


Fig. 7. Visualization of attention maps from Baseline, MHSAB and MHSAB+ACM. As shown in column four and eight, the attention areas from MHSAB+ACM can effectively locate on the key suregions.

Based on the MHSAB, we introduce the ACM into MHSAB, i.e. MHSAB+ACM. Compared with MHSAB, some attention information is suppressed and some attention information is highlighted. And subjectively, we can see that the highlight attention areas conducted by MHSAB+ACM are more important than that of MHSAB. Besides, the left half of the Fig. 7, compared with baseline model, our MHSAB, and

MHSAB+ACM can better capture the key information from the unoccluded regions in the occlusion person images.

Besides, as shown in Fig. 8, we visualized the feature maps of the 8 feature branches (from K_1 to K_8) of MHSAB under different variants. The feature maps corresponding to K_i , ($i = 1, 2, \dots, K$) represents the person information captured by the K_i -th attention head in MHSAB. (I) As shown in “Group A: MHSAB” of Fig. 8, the Feature Diversity Regularization Term (FDRT) can help the MHSAB effectively capture diversity information for person matching. Since this, MHSAM can capture key information from different perspectives, which enhances the model robustness. (II) If we remove the FDRT from MHSAB, i.e. “Group B: MHSAB w/o FDRT”, the feature maps (from K_1 to K_8) are mixed with some redundant information, and some feature responses (in red box especially) are scattered and weak. (III) If we remove the FDRT and Improved Hard Triplet Loss(IHTL) from MHSAB, i.e. “Group B: MHSAB w/o FDRT and IHTL”, the responses of key information in the feature maps become sparse and weak. Without the constraint of IHTL and FDRT, it is difficult to ensure that every feature branch head in MHSAB can capture the key diversity information for person matching.

In all, results in Table V can sufficiently evidence that the effectiveness of each component.

E. Hyper-parameters Discussion in MHSA-Net

1) Multi-Head Self-Attention Mechanism (MHSAM):

In this module, we (1) try to find suitable K for MHSAM; and (2) discuss the influence of different mechanisms in MHSAM.

(1) Table VII shows the main results in the Market-1501 and CUHK03-Detected datasets. We set $\lambda_1 = 0$ and $\lambda_2 = 0$ in Eq. 11, and $\lambda_3 = 0$ in Eq. 13. Here, we set $K = 1, 2, \dots, 9$ in MHSAM. As the results show, MHSAM with $K > 0$ can effectively improve the baseline’s performance in both datasets. And we find $K = 7, 8$ are suitable parameters when MHSAM produces good performance in both datasets over both measures.

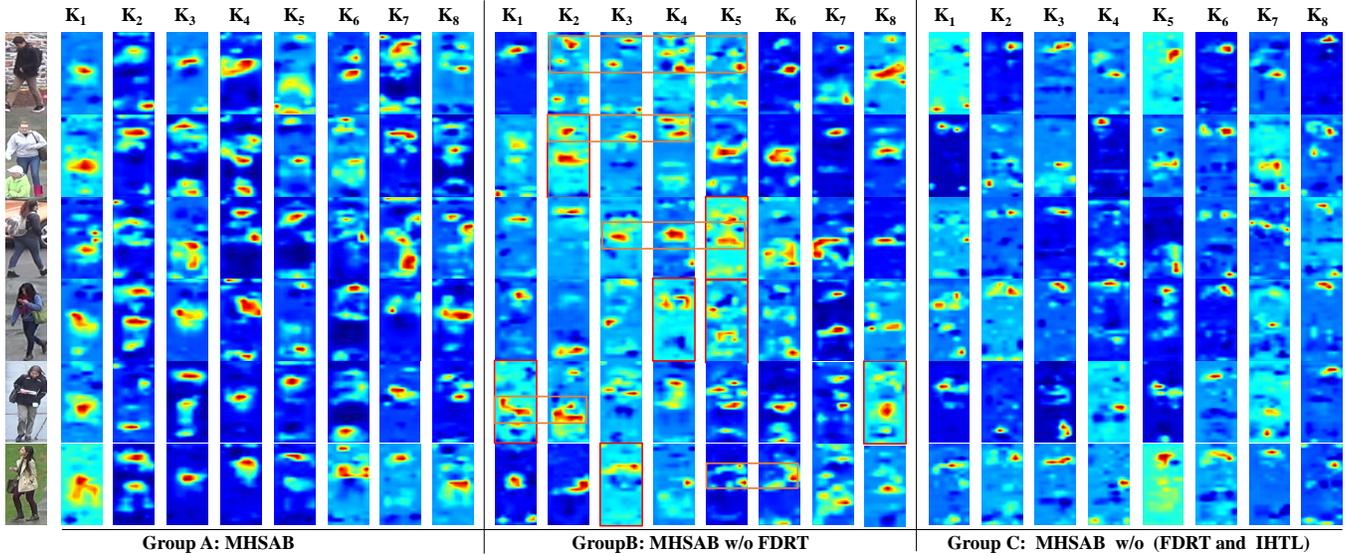


Fig. 8. The feature maps of the $K = 8$ feature branches (from K_1 to K_8) of MHSAB under different variants. The feature maps corresponding to K_i , ($i = 1, 2, \dots, K$) represents the person information captured by the K_i -th attention head in MHSAB.

TABLE VII

RESULTS ON MARKET-1501 AND CUHK03-DETECTED TESTING DATA FROM DIFFERENT HYPER-PARAMETER K IN THE MHSAM. THE BOLD IS THE BEST RESULT.

Method	Market-1501		CUHK03-Detected	
	Rank-1	mAP	Rank-1	mAP
Baseline	92.0	78.8	56.3	53.0
$K = 1$	93.0	80.8	65.5	61.2
$K = 2$	93.0	80.8	66.1	62.0
$K = 3$	93.1	80.7	66.1	62.9
$K = 4$	93.2	80.5	68.3	64.0
$K = 5$	92.7	80.6	68.2	64.8
$K = 6$	92.8	80.0	67.4	64.6
$K = 7$	93.3	81.0	68.9	65.1
$K = 8$	93.0	80.2	68.1	64.4
$K = 9$	92.6	80.1	68.6	65.4

TABLE VIII

INFLUENCE OF DIFFERENT STRATEGY ON THE OCCLUDED-DUKEMTMC DATASET IN THE MHSAM. HERE, THE HYPER-PARAMETER K IN THE MHSAM IS 8. THE BOLD IS THE BEST RESULT.

Method	Occluded-DukeMTMC			
	Rank-1	Rank-5	Rank-10	mAP
MHSA-Net*	59.7	74.3	79.5	44.8
MHSA-Net* (CONCAT)	51.7	68.1	73.5	36.3
MHSA-Net* (SUM)	53.8	70.9	76.4	38.8
MHSA-Net* w/o RLM	55.9	72.4	77.4	41.8

TABLE IX

RESULTS ON CUHK03-DETECTED TESTING DATA FROM DIFFERENT HYPER-PARAMETERS λ_1 IN FEATURE DIVERSITY REGULARIZATION TERM (FDRT). THE BOLD IS THE BEST RESULT.

Method	CUHK03-Detected	
	Rank-1	mAP
Baseline+MHSAM	68.1	64.4
$K = 8, \lambda_1 = 10^{-6}$	69.6	65.8
$K = 8, \lambda_1 = 10^{-5}$	69.7	64.9
$K = 8, \lambda_1 = 10^{-4}$	70.0	65.7
$K = 8, \lambda_1 = 10^{-3}$	69.4	65.3
$K = 8, \lambda_1 = 10^{-2}$	69.1	65.4
$K = 8, \lambda_1 = 10^{-1}$	69.1	65.3

(2) Table VIII show the influence of different mechanisms in MHSAM on the Occluded-DukeMTMC dataset. First, the effects of different feature fusion operations to $P^\perp(x) \in \mathbb{R}^{K \times 512}$ in MHSAM are compared: MHSA-Net* uses SAFFM to fuse the $P^\perp(x)$ and produces a 512-dimensional vector; MHSA-Net* (CONCAT) directly concatenates the $P^\perp(x)$ to one $(K + 1) \times 512$ vector; MHSA-Net* (SUM) simply sums up the $P^\perp(x)$ to one 512-D vector. The SAFFM results in the best performance in these three feature fusion operations. The output vector from SAFFM also has much lower dimension than that from feature concatenation. Second, the effect of the “Residual Learning Module” is compared. “MHSA-Net* w/o RLM” drops the “Residual Learning Module” from MHSA-Net*, and this leads to declined performance. Hence, “Residual Learning Module” is beneficial for MHSA-Net* to capture useful local information.

2) Feature Diversity Regularization Term (FDRT):

This section discusses the suitable λ_1 for \mathcal{L}_{FDRT} , and FDRT’s performance under different values of hyper-parameter K in MHSAM. Table IX and Fig. 9 show the

results. We set $\lambda_2 = 0$ (Eq. 11) and $\lambda_3 = 0$ (Eq. 13).

Table IX shows Rank-1 and mAP results under different λ_1 , from 10^{-6} to 10^{-1} . Both Rank-1 and mAP reaches the highest score when $\lambda_1 = 10^{-4}$.

Then, with this setting of $\lambda_1 = 10^{-4}$ (Eq. 6), we discuss performance of FDRT on different hyper-parameter K values in MHSAM. We conducted ablation studies in Market-1501 and CUHK03-Detected datasets. Compared with MHSAM (in Orange polylines and bars Figs. 9), adding FDRT in MHSAM (in Blue polylines and bars, respectively) improves the performance of MHSAM in person Re-ID.

3) Improved Hard Triplet Loss (IHTL):

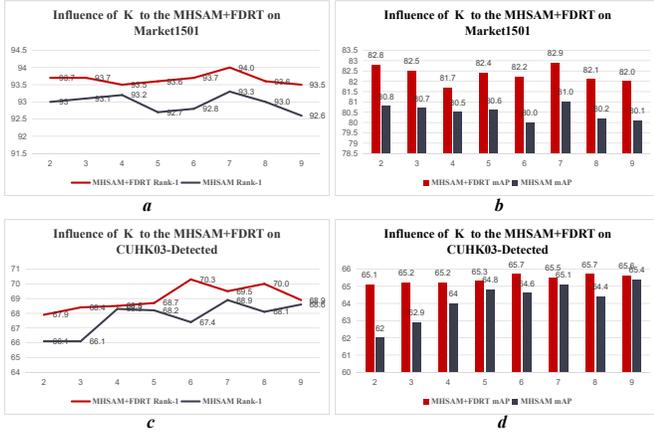


Fig. 9. Results of MHSAM+FDRT and MHSAM in the Market-1501 and CUHK03-Detected testing data from different hyper-parameters K in MHSAM. Here, the hyper-parameter $\lambda_1 = 10^{-4}$. Figure a shows the influence of the parameter K on the MHSAM+FDRT under Rank-1 score in Market-1501 dataset. Figure b shows the influence of the parameter K on the MHSAM+FDRT under mAP score in Market-1501 dataset. Figure c shows the influence of the parameter K on the MHSAM+FDRT under Rank-1 score in CUHK03-Detected dataset. Figure d shows the influence of the parameter K on the MHSAM+FDRT under mAP score in CUHK03-Detected dataset.

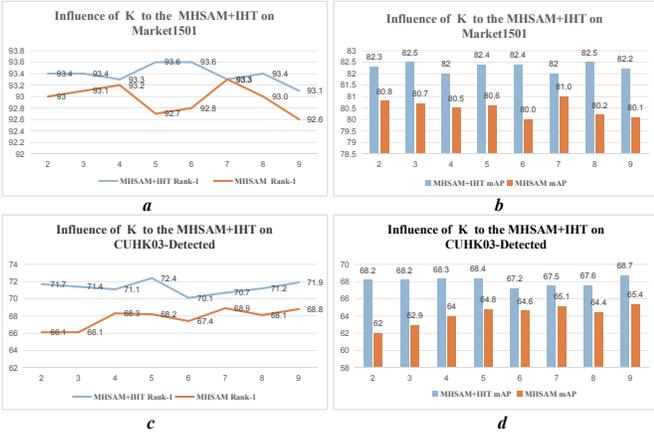


Fig. 10. Results of the MHSAM+IHTL and the MHSAM in the Market-1501 and CUHK03-Detected testing data from different hyper-parameters K in MHSAM. Here, the hyper-parameter $\lambda_2 = 1.0$. Figure a shows the influence of the parameter K on the MHSAM+IHTL under Rank-1 score in Market-1501 dataset. Figure b shows the influence of the parameter K on the MHSAM+IHTL under mAP score in Market-1501 dataset. Figure c shows the influence of the parameter K on the MHSAM+IHTL under Rank-1 score in CUHK03-Detected dataset. Figure d shows the influence of the parameter K on the MHSAM+IHTL under mAP score in CUHK03-Detected dataset.

For this module, we find the suitable weight λ_2 in IHTL, and discuss the performance of different K values in MHSAM. Table X and Fig. 10 show the main results. Here we set $\lambda_1 = 0$ (Eq. 11) and $\lambda_3 = 0$ (Eq. 13).

Table X shows that with the setting of λ_2 from 0.01 to 100, Rank-1 and mAP scores get better and then decline. For both $K = 7$ and 8, IHTL gets the best performance when $\lambda_2 = 1.0$.

Then, with $\lambda_2 = 1.0$ set, we compare the performance of different K in MHSAM. We conducted ablation studies in Market-1501 and CUHK03-Detected datasets. Compared with the MHSAM (Orange polylines and bars in Fig. 10),

TABLE X
RESULTS ON CUHK03-DETECTED TESTING DATA FROM DIFFERENT HYPERPARAMETERS λ_2 IN IMPROVED HARD TRIPLET LOSS (IHTL). HERE $K = 7, 8$ IN THE MHSAM. THE BOLD IS THE BEST RESULT.

Method	CUHK03-Detected	
	Rank-1	mAP
$K = 7$ Baseline+MHSAM	68.9	65.1
$K = 7, \lambda_2 = 10^2$	61.2	59.5
$K = 7, \lambda_2 = 10^1$	67.8	65.5
$K = 7, \lambda_2 = 1.0$	70.7	67.5
$K = 7, \lambda_2 = 10^{-1}$	69.8	65.2
$K = 7, \lambda_2 = 10^{-2}$	68.6	64.8
$K = 8$ Baseline+MHSAM	68.1	64.4
$K = 8, \lambda_2 = 10^2$	63.4	61.1
$K = 8, \lambda_2 = 10^1$	69.9	67.0
$K = 8, \lambda_2 = 1.0$	71.2	67.6
$K = 8, \lambda_2 = 10^{-1}$	70.9	67.5
$K = 8, \lambda_2 = 10^{-2}$	68.7	65.1

introducing IHTL into MHSAM, i.e. MHSAM+IHTL (Blue polylines and bars in Fig. 10), improves the performance of MHSAM in person Re-ID.

F. Attention Competition Mechanism (ACM)

In this module, we try to find suitable hyper-parameters λ_3 (Eq. 13) and γ (Eq. 12) in the Attention Competition Mechanism (ACM). Table XI shows the main results. We conducted the experiments on the CUHK03-Detected test dataset with $K = 8$ set in MHSAM.

The competition loss functions in ACM are $\mathcal{L}_{MHSAB} = \mathcal{L}_{SAFFM} + \lambda_1 \mathcal{L}_{FDRT} + \lambda_2 \mathcal{L}_{IHTL} + \mathcal{L}_{ReN}$ and \mathcal{L}_C . The \mathcal{L}_{FDRT} and \mathcal{L}_{IHTL} are two terms we proposed here. To demonstrate the effectiveness of ACM individually, we set $\lambda_1 = 0$ and $\lambda_2 = 0$ in \mathcal{L}_{MHSAB} . (Table V in Section VII-D shows the effectiveness of ACM combined with other contributions.)

Table XI shows that when $\gamma \leq 10^{-2}$ and $\lambda_3 \leq 10^{-1}$, the ACM effectively improves the performance of MHSAM in the person Re-ID task. When γ is too big, the performance declines. This is because that each element in α belongs to $[0, 1]$. If we set a too-big γ , attention weights on most regions will be suppressed by \mathcal{L}_C . Based on these observations, we set $\lambda_3 = 10^{-3}$ and $\gamma = 10^{-3}$ in our MHSAM-Net. In all, based on the suitable parameter setting, ACM can effectively improve the performance of the MHSAB.

VIII. LIMITATION AND DISCUSSION

Although our proposed MHSAM-Net achieved the competitive performance in the occluded person Re-ID task, some limitations and discussion must be taken into consideration.

Firstly, our proposed MHSAM-Net mainly considers the situation of objects occluding person. The situation of person occluding person is not considered in the proposed MHSAM-Net. Therefore, it is necessary to continue to optimize MHSAM-Net in future work, to improve the model performance on the situation of person occluding person.

Besides, the person search requires a combination of the person detection task and the person Re-ID task in actual

TABLE XI

RESULTS ON CUHK03-DETECTED TESTING DATA FROM THE DIFFERENT HYPERPARAMETER λ_3 AND γ IN ATTENTION REGULARIZATION TERM 7. HERE, $K = 8$ IN THE MHSAM. THE BOLD IS THE BEST RESULT.

Method	CUHK03-Detected	
	Rank-1	mAP
Baseline+MHSAM	68.1	64.4
$\lambda_3 = 10^{-5}, \gamma = 10^{-3}$	68.7	65.1
$\lambda_3 = 10^{-4}, \gamma = 10^{-3}$	69.1	65.0
$\lambda_3 = 10^{-3}, \gamma = 10^{-3}$	69.9	65.8
$\lambda_3 = 10^{-2}, \gamma = 10^{-3}$	69.0	65.3
$\lambda_3 = 10^{-1}, \gamma = 10^{-3}$	69.4	65.1
$\lambda_3 = 1, \gamma = 10^{-3}$	58.4	55.6
$\lambda_3 = 10^{-3}, \gamma = 10^{-4}$	68.8	65.9
$\lambda_3 = 10^{-3}, \gamma = 10^{-2}$	69.5	65.4
$\lambda_3 = 10^{-3}, \gamma = 10^{-1}$	67.9	64.5
$\lambda_3 = 10^{-3}, \gamma = 5 \times 10^{-1}$	67.4	64.2
$\lambda_3 = 10^{-3}, \gamma = 1$	66.9	63.8

scenarios. Only relying on the person Re-ID model cannot effectively search for the target person. Therefore, it is necessary to study how to combine MHSAN-Net with person detection models to build an end-to-end person search framework

IX. CONCLUSION

We proposed a Multi-Head Self-Attention Network (MHSA-Net) to improve the ability of the Re-ID model on capturing the key information from the occluded person image. Specifically, we introduced the Multi-Head Self-Attention Branch (MHSAB) to adaptively capture key local person information, and produce multiple diversity embeddings of one person image to facilitate person matching. We also designed an Attention Competition Mechanism (ACM) to further help MHSAB prune out non-important local information. Extensive experiments were conducted to validate the effectiveness of each component in MHSA-Net; and they showed that MHSA-Net achieves competitive performance on three standard person Re-ID datasets and four occlusion person Re-ID datasets.

ACKNOWLEDGMENTS

This work is supported by National Key R&D Program of China (No. 2021ZD0111900), National Natural Science Foundation of China (No.U21B2038, 61876012, 61976040, 62172073, 61906011), National Science Foundation of USA (OIA-1946231, CBET-2115405), Chinese Postdoctoral Science Foundation (No. 2021M700303), Liaoning Provincial Natural Science Foundation of China (No. 2021-MS-110). No conflict of interest: Hongchen Tan, Xiuping Liu, Baocai Yin, and Xin Li declare that they have no conflict of interest.

REFERENCES

[1] Hermans Alexander, Beyer Lucas, and Leibe. Bastian. In defense of the triplet loss for person re-identification. In *arXiv:1703.07737*, 2017.

[2] He Anfeng, Luo Chong, Tian Xinmei, and Zeng Wenjun. A twofold siamese network for real-time object tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[3] Vaswani Ashish, Shazeer Noam, Parmar Niki, Uszkoreit Jakob, Jones Llion, N Gomez n Łukasz Kaiser Aidan, and Polosukhin Illia. A structured self-attentive sentence embedding. In *Conference and Workshop on Neural Information Processing Systems*, 2017.

[4] S. Bai, X. Bai, and Q. Tian. Scalable person re-identification on supervised smoothed manifold. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3356–3365, Los Alamitos, CA, USA, jul 2017. IEEE Computer Society.

[5] Xia Bryan (Ning), Gong Yuan, Zhang Yizhe, and Poellabauer Christian. Second-order non-local attention networks for person re-identification. In *IEEE International Conference on Computer Vision*, 2019.

[6] Xuesong Chen, Canmiao Fu, Yong Zhao, Feng Zheng, Jingkuan Song, Rongrong Ji, and Yi Yang. Saliency-guided cascaded suppression network for person re-identification. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3297–3307, 2020.

[7] Wang Cheng, Zhang Qian, Huang Chang, Liu Wenyu, and Wang Xinggong. Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[8] Su Chi, Li Jianing, Zhang Shiliang, Xing Junliang, Gao Wen, and Tian Qi. Pose-driven deep convolutional model for person re-identification. In *IEEE International Conference on Computer Vision*, 2017.

[9] Tay Chiat-Pin, Roy Sharmili, and Yap Kim-Hui. Aanet: Attribute attention network for person re-identifications. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[10] Song Chunfeng, Huang Yan, Wanli Ouyang, and Liang Wang. Mask-guided contrastive attention model for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1179–1188, 2018.

[11] Lin Chunli and Wang Kejun. A behavior classification based on enhanced gait energy image. In *2010 International Conference on Networking and Digital Society*, volume 2, pages 589–592, 2010.

[12] Yongxing Dai, Jun Liu, Yan Bai, Zekun Tong, and Ling-Yu Duan. Dual-refinement: Joint label and feature refinement for unsupervised domain adaptive person re-identification. *IEEE Transactions on Image Processing*, 30:7815–7829, 2021.

[13] Xie Di, Xiong Jiang, and Pu Shiliang. All you need is beyond a good init: Exploring better solution for training extremely deep convolutional neural networks with orthonormality and modulation. In *CVPR*, 2017.

[14] Jian Dong, Qiang Chen, Shen Xiaohui, Yang Jianchao, and Yan Shuicheng. Towards unified human parsing and pose estimation. In *CVPR*, 2014.

[15] Bahdanau Dzmitry, Cho Kyunghyun, and Bengio Yoshua. Neural machine translation by jointly learning to align and translate. In <https://arxiv.org/abs/1409.0473v2>.

[16] Liu Fangyi and Zhang Lei. View confusion feature learning for person re-identification. In *IEEE International Conference on Computer Vision*, 2019.

[17] Shang Gao, Jingya Wang, Huchuan Lu, and Zimo Liu. Pose-guided visible part matching for occluded person reid. In *CVPR*, 2020.

[18] Wang Guan'an, Yang Shuo, Liu Huanyu, Wang Zhicheng, Yang Yang, Wang Shuliang, Yu Gang, Zhou Erjin, and Sun Jian. High-order information matters: Learning relation and topology for occluded person re-identification. In *CVPR*, 2020.

[19] Wang Guanshuo, Yuan Yufeng, Chen Xiong, Li Jiwei, and Zhou Xi. Learning discriminative features with multiple granularities for person re-identification. In *arXiv:1804.01438*, 2018.

[20] Albert Haque, Alexandre Alahi, and Li Fei-Fei. Recurrent attention models for depth-based person identification. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1229–1238, 2016.

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[22] Tan Hongchen, Liu Xiuping, Li Xin, Zhang Yi, and Yin Baocai. Semantics-enhanced adversarial nets for text-to-image synthesis. In *ICCV*, 2019.

[23] Xu Hongyu, Wang Zhangyang, Yang Haichuan, Liu Ding, and Liu Ji. Learning simple thresholded features with sparse support recovery. In *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.

[24] Huang Houjing, Li Dangwei, Zhang Zhang, Chen Xiaotang, and Huang Kaiqi. Adversarially occluded samples for person re-identification. In *CVPR*, 2018.

[25] Cordonnier Jean-Baptiste, Loukas Andreas, and Jaggi Martin. On the relationship between self-attention and convolutional layers. In *ICLR*, 2020.

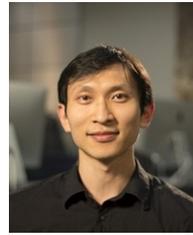
[26] Deng Jia, Dong Wei, Socher Richard, Li Li-Jia, Li Kai, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.

[27] Si Jianlou, Zhang Honggang, Li Chun-Guang, Kuen Jason, Kong Xi-angfei, Kot Alex C, and Wang Gang. Dual attention matching network for context-aware feature sequence based person re-identification. In *arXiv preprint arXiv:1803.09937*, 2018.

[28] Lu Jiasen, Xiong Caiming, Parikh Devi, and Socher Richard. Knowing

- when to look: Adaptive attention via a visual sentinel for image captioning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [29] Miao Jiayu, Wu Yu, Liu Ping, Ding Yuhang, and Yang Yi. Pose-guided feature alignment for occluded person re-identification. In *IEEE International Conference on Computer Vision*, 2019.
- [30] Zhuo Jiaxuan, Chen Zeyu, Lai Jianhuang, and Wang Guangcong. Occluded person re-identification. In *ICME*, 2018.
- [31] Ba Jimmy Lei, Kiros Jamie Ryan, and Hinton Geoffrey E. Layer normalization. In *arXiv preprint arXiv:1607.06450*, 2016.
- [32] Xu Jing, Zhao Rui, Zhu Feng, Wang Huaming, and Ouyang Wanli. Attention-aware compositional network for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [33] Mahdi M. Kalayeh, Emrah Basaran, Muhittin Gokmen, Mustafa E. Kamasak, and Mubarak Shah. Human semantic parsing for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [34] M. Esat Kalfaoglu, Sinan Kalkan, and A. Aydin Alatan. Late temporal modeling in 3d CNN architectures with BERT for action recognition. *CoRR*, abs/2008.01232, 2020.
- [35] Nikolaos Karianakis, Zicheng Liu, Yinpeng Chen, and Stefano Soatto. Person depth reid: Robust person re-identification with commodity depth sensors. *CoRR*, abs/1705.09882, 2017.
- [36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [37] Jun Li, Xianglong Liu, Wenxuan Zhang, Mingyuan Zhang, Jingkuan Song, and Nicu Sebe. Spatio-temporal attention networks for action recognition and detection. *IEEE Transactions on Multimedia*, 22(11):2990–3001, 2020.
- [38] Zheng Liang, Shen Liyue, Tian Lu, Wang Shengjin, Wang Jingdong, and Tian Qi. Scalable person re-identification: A benchmark. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1116–1124, 2015.
- [39] He Lingxiao, Liang Jian, Li Haiqing, and Sun Zhenan. Deep spatial feature reconstruction for partial person reidentification: Alignment-free approach. In *CVPR*, 2018.
- [40] He Lingxiao, Wang Yinggang, Liu Wu, Liao Xingyu, Zhao He, Sun Zhenan, and Feng Jiashi. Foreground-aware pyramid reconstruction for alignment-free occluded person re-identification. In *ICCV*, 2019.
- [41] He Lingxiao, Sun Zhenan, Zhu Yuhao, and Wang Yunbo. Recognizing partial biometric patterns. In *arXiv preprint arXiv:1810.07399*, 2018.
- [42] Matteo Munaro, Andrea Fossati, Alberto Basso, Emanuele Menegatti, and Luc Van Gool. One-shot person re-identification with a consumer depth camera. In Shaogang Gong, Marco Cristani, Shuicheng Yan, and Chen Change Loy, editors, *Person Re-Identification*, Advances in Computer Vision and Pattern Recognition, pages 161–181. Springer, 2014.
- [43] Xin Ning, Ke Gong, Weijun Li, and Liping Zhang. Jwsaa: Joint weak saliency and attention aware for person re-identification. *Neurocomputing*, 453:801–811, 2021.
- [44] Goyal Priya, Dollar Piotr, Girshick Ross, Noord-huis Pieter, Wesolowski Lukasz, Kyrola Aapo, Tulloch Andrew, Jia Yangqing, and He. Kaiming. Accurate, large minibatch sgd: Training imagenet in 1 hour. In *arXiv:1706.02677*, 2017.
- [45] Haocong Rao, Siqi Wang, Xiping Hu, Mingkui Tan, Yi Guo, Jun Cheng, Xinwang Liu, and Bin Hu. A self-supervised gait encoding approach with locality-awareness for 3d skeleton based person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [46] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *European Conference on Computer Vision (ECCV)*, 2016.
- [47] Hou Ruibing, Ma Bingpeng, Chang Hong, Gu Xinqian, Shan Shiguang, and Chen Xilin. Interaction-and-aggregation network for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [48] Zhou Sanping, Wang Jinjun, Meng Deyu, Liang Yudong, Gong Yihong, and Zheng Nanning. Discriminative feature learning with foreground attention for person re-identification. *IEEE Transactions on Image Processing*, 28(9):4671 – 4684, 2019.
- [49] Liao Shengcai, Jain Anil K, and Li Stan Z. Partial face recognition: Alignment-free approach. In *TPAMI*, 2013.
- [50] Sabesan Sivapalan, Daniel Chen, Simon Denman, Sridha Sridharan, and Clinton Fookes. Gait energy volumes and frontal gait recognition using depth images. In *2011 International Joint Conference on Biometrics (IJCB)*, pages 1–6, 2011.
- [51] Chen Tianlong, Ding Shaojin, Xie Jingyi, Yuan Ye, Chen Wuyang, Yang Yang, Ren Zhou, and Wang Zhangyang. Abd-net: Attentive but diverse person re-identification. In *IEEE International Conference on Computer Vision*, 2019.
- [52] Zheng Wang, Junjun Jiang, Yang Wu, Mang Ye, Xiang Bai, and Shinichi Satoh. Learning sparse and identity-preserved hidden attributes for person re-identification. *IEEE Transactions on Image Processing*, 29:2013 – 2025, 2020.
- [53] Li Wei, Zhao Rui, Xiao Tong, and Wang Xiaogang. Deepreid: Deep filter pairing neural network for person reidentification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [54] Li Wei, Zhu Xiayan, and Gong Shaogang. Harmonious attention network for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2285–2294, 2018.
- [55] Zheng Wei-Shi, Li Xiang, Xiang Tao, Liao Shengcai, Lai Jianhuang, and Gong Shaogang. Partial person reidentification. In *ICCV*, 2015.
- [56] Yang Wenjie, Huang Houjing, Zhang Zhang, Chen Xiaotang, Huang Kaiqi, and Zhang. Shu. Towards rich feature discovery with class activation maps augmentation for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [57] Guile Wu, Xiayan Zhu, and Shaogang Gong. Learning hybrid ranking representation for person re-identification. *Pattern Recognition*, 121:108239, 2022.
- [58] Wanyin Wu, Dapeng Tao, Hao Li, Zhao Yang, and Jun Cheng. Deep features for person re-identification on metric learning. *Pattern Recognition*, 110:107424, 2021.
- [59] Chang Xiaobin, M Hospedales Timothy, and Xiang Tao. Multi-level factorisation net for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [60] Qian Xuelin, Fu Yanwei, Xiang Tao, Wang Wenxuan, Qiu Jie, Wu Yang, Jiang Yu-Gang, and Xue Xiangyang. Pose normalized image generation for person re-identification. In *European Conference on Computer Vision*, 2018.
- [61] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven C.H. Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [62] Cho Yeong-Jun and Yoon Kuk-Jin. Improving person reidentification via pose-aware multi-shot matching. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [63] Sun Yifan, Zheng Liang, Yang Yi, Tian Qi, and Wang Shengjin. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [64] Ge Yixiao, Li Zhuowan, Zhao Haiyu, Yin Guojun, Yi Shuai, and Wang Xiaogang. Fd-gan: Pose-guided feature distilling gan for robust person re-identification. In *NIPS*, 2018.
- [65] Suh Yumin, Wang Jingdong, Tang Siyu, Mei Tao, and Lee Kyoung Mu. Part-aligned bilinear representations for person re-identification. In *ECCV*, 2018.
- [66] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. Relation-aware global attention for person re-identification. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3183–3192, 2020.
- [67] Bin Zhao, Maoguo Gong, and Xuelong Li. Hierarchical multimodal transformer to summarize videos. *Neurocomputing*, 468:360–369, 2022.
- [68] Bin Zhao, Haopeng Li, Xiaoqiang Lu, and Xuelong Li. Reconstructive sequence-graph network for video summarization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [69] Bin Zhao, Xuelong Li, and Xiaoqiang Lu. Cam-rnn: Co-attention model based rnn for video captioning. *IEEE Transactions on Image Processing*, 28(11):5552–5565, 2019.
- [70] Cairong Zhao, Kang Chen, Zhihua Wei, Yipeng Chen, Duoqian Miao, and Wei Wang. Multilevel triplet deep learning model for person re-identification. *Pattern Recognition Letters*, 117:161–168, 2019.
- [71] Cairong Zhao, Xinbi Lv, Shuguang Dou, Shanshan Zhang, Jun Wu, and Liang Wang. Incremental generative occlusion adversarial suppression network for person reid. *IEEE Transactions on Image Processing*, 30:4212–4224, 2021.
- [72] Cairong Zhao, Xinbi Lv, Zhang Zhang, Wangmeng Zuo, Jun Wu, and Duoqian Miao. Deep fusion feature representation learning with hard mining center-triplet loss for person re-identification. *IEEE Transactions on Multimedia*, 22(12):3180–3195, 2020.
- [73] Cairong Zhao, Xuekuan Wang, Wangmeng Zuo, Fumin Shen, Ling Shao, and Duoqian Miao. Similarity learning with joint transfer constraints for person re-identification. *Pattern Recognition*, 97:107014, 2020.
- [74] Zheng Zhedong, Zheng Liang, and Yang Yi. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *International Conference on Computer Vision (ICCV)*, 2017.
- [75] Zheng Zhedong, Zheng Liang, and Yang Yi. Pedestrian alignment network for large-scale person re-identification. *IEEE Transactions on*

- Circuits and Systems for Video Technology*, 29(10):3037–3045, 2019.
- [76] Zhedong Zheng and Yi Yang. Person re-identification in the 3d space. *CoRR*, abs/2006.04569, 2020.
- [77] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Learning generalisable omni-scale representations for person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [78] Lin Zhouhan, Feng Minwei, Nogueira dos Santos Cicero, Yu Mo, Xiang Bing, Zhou Bowen, and Bengio Yoshua. A structured self-attentive sentence embedding. In *International Conference on Learning Representations*, 2017.
- [79] Zhong Zhun, Zheng Liang, Cao Donglin, and Li Shaozi. Re-ranking person re-identification with k-reciprocal encoding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [80] Zhong Zhun, Zheng Liang, Kang Guoliang, Li Shaozi, and Yang Yi. Random erasing data augmentation. In *arXiv:1708.04896*, 2017.
- [81] Zhong Zhun, Zheng Liang, Kang Guoliang, Li Shaozi, and Yang Yi. Random erasing data augmentation. In *The National Conference on Artificial Intelligence*, 2020.
- [82] Yang Zichao, He Xiaodong, Gao Jianfeng, Deng Li, and Smola Alex. Stacked attention networks for image question answering. In *CVPR*, 2016.
- [83] Dai Zuozhuo, Chen Mingqiang, Gu Xiaodong, Zhu Siyu, and Tan Ping. Batch dropout network for person re-identification and beyond. In *IEEE International Conference on Computer Vision*, 2019.



Xin Li is a Professor at Division of Electrical & Computer Engineering, Louisiana State University, USA. He got his B.E. degree in Computer Science from University of Science and Technology of China in 2003, and his M.S. and Ph.D. degrees in Computer Science from State University of New York at Stony Brook in 2005 and 2008. His research interests are in Geometric and Visual Data Computing, Processing, and Understanding, Computer Vision, and Virtual Reality. For more detail, please see <http://www.ece.lsu.edu/xinli>.



Hongchen Tan is a Lecturer of Artificial Intelligence Research Institute at Beijing University of Technology. He received Ph.D degrees in computational mathematics from the Dalian University of Technology in 2021. His research interests are person Re-identification, Image Synthesis, and Object Detection. Various parts of his work have been published in top conferences and journals, such as IEEE TIP/TMM/TCSVT/TNNLS/ICCV, and Neuro-computing.



Xiuping Liu is a Professor in School of Mathematical Sciences at Dalian University of Technology. She received Ph.D degrees in computational mathematics from Dalian University of Technology in 1999. Her research interests include shape modeling and analyzing, and computer vision.



Baocai Yin is a Professor of Artificial Intelligence Research Institute at Beijing University of Technology. He is also a Researcher with the Beijing Key Laboratory of Multimedia and Intelligent Software Technology and the Beijing Advanced Innovation Center for Future Internet Technology. He received the M.S. and Ph.D. degrees in computational mathematics from the Dalian University of Technology, Dalian, China, in 1988 and 1993, respectively. His research interests include multimedia, image processing, computer vision, and pattern recognition.