

Part-Attention Based Model Make Occluded Person Re-Identification Stronger

1st Zhihao Chen

Computer School

Beijing Information Science and Technology University

Beijing, China

2021011561@bistu.edu.cn

2nd Yiyuan Ge*

School of Instrument Science and Opto-Electronics Engineering

Beijing Information Science and Technology University

Beijing, China

geiyuan@bistu.edu.cn

Abstract—The goal of occluded person re-identification (ReID) is to retrieve specific pedestrians in occluded situations. However, occluded person ReID still suffers from background clutter and low-quality local feature representations, which limits model performance. In our research, we introduce a new framework called PAB-ReID, which is a novel ReID model incorporating part-attention mechanisms to tackle the aforementioned issues effectively. Firstly, we introduce the human parsing label to guide the generation of more accurate human part attention maps. In addition, we propose a fine-grained feature fuser for generating fine-grained human local feature representations while suppressing background interference. Moreover, We also design a part triplet loss to supervise the learning of human local features, which optimizes intra/inter-class distance. We conducted extensive experiments on specialized occlusion and regular ReID datasets, showcasing that our approach outperforms the existing state-of-the-art methods.

Index Terms—Occluded ReID, attention maps, human parsing labels.

I. INTRODUCTION

Person re-identification (ReID) aims to retrieve specific pedestrians from different scenes and camera views and is of interest due to its wide range of applications in security [1]. Existing ReID methods learn by extracting a global representation of the target person. As shown in Figure 1(a), occlusion impacts the visual representation of pedestrians. When occlusion occurs, the distinction between various categories diminishes, leading to a situation where images with different IDs may have similar global representations. Also, as shown in Fig. 1(b), various occlusions widen the distance within the same category, which suggests that the global representations of the same pedestrian may be different. In summary, different parts of the occlusion and occluder contents can easily lead to wrong detection results.

There are two primary mainstream approaches for solving occluded ReID [50],[51],[52],[53],[54],[55],[56],[57],[58],[59],[60],[61],[62],[63],[64],[65],[66],[67],[68],[69],[77],[78], namely extra information-based methods and part-to-part matching methods. Among them, extra information-based methods use only the visible part of the body for ReID [2–5], which are usually based on extra information to judge the visibility of



Fig. 1. Examples of Challenges in Occluded Person Re-identification. In Figure. 1(a), similar occlusions reduce the gap between different categories. In Figure. 1(b), different occlusions can lead to an increase in distances within the same category. Figure. 1(c) illustrates that the same body part may have a similar appearance across different individuals.

the body part (e.g., pose estimation). The above methods are computationally expensive and not robust when faced with complex occlusion situations (e.g., when occluding between pedestrians). The majority of the current approaches are based on part-to-part matching, which solves the occlusion problem by comparing the similarity between local features [38–40]. However, this part-to-part matching approach still faces some challenges: 1) The part-to-part matching-based approach relies on spatial attention maps [30] to construct body part features for ReID targets. However, the current ReID dataset lacks explicit annotations about the body part regions. 2) In occlusion scenes, the challenge remains to suppress background clutter and extract fine-grained effective features. 3) As shown in Fig. 1(c), the same body part may have a similar appearance across individuals. In other words, the visual appearance of body parts may not be inherently distinctive, making the conventional ReID loss, designed for learning global representations, less effective when applied to part representation learning.

To solve the above problems, we propose the following solutions. **Firstly**, we introduce human parsing labels to guide the generation of attention maps and use pixel-level attention predictors to predict the attribution of each pixel to generate detailed attention maps. At the same time, dual-loss supervised training is used, with both body part prediction targets

*Corresponding author

and ReID targets. This dual-supervision mechanism makes the final obtained attention map more relevant to the ReID task. **Secondly**, we then use a fine-grained feature focuser to process the generated attention maps. The effect of the attention map is enhanced by filtering irrelevant background information, and finally, fine-grained body part ReID features are generated. **Thirdly**, the standard ReID loss function assumes that different individuals have different appearances, i.e., different global feature vectors. The above assumption does not hold when confronted with part-based feature vectors. For this reason, we propose part triplet loss to supervise learning, which enhances the robustness of the model to similar part features.

Finally, we combine the part attention block with the global-local learning block to propose the part-attention based model PAB-ReID. Here are our contributions:

- We propose a part-attention based ReID model (PAB-ReID), which achieves person re-identification in the occlusion case by learning attention maps of different body parts.
- We innovatively introduce the human parsing labels to guide the generation of the attention maps, which leads to more precise feature extraction regions for each body part.
- We design a fine-grained feature focuser in a global-local learning block for part features, which can filter irrelevant background information and generate fine-grained body part features.
- Part triplet loss is proposed for supervised learning of body part features, which is robust to similar body part features.

II. RELATED WORKS

A. Person Re-Identification

Person re-identification identifies and matches target pedestrians from existing video sequences, and most of the mainstream methods use CNN [70–76] architecture. The widespread effectiveness of convolutional neural networks in image processing has resulted in a growing adoption of neural network-based approaches for addressing ReID tasks, and the mainstream approaches can be divided into three categories. The first method uses local features to capture locally salient information about pedestrians by dividing the output of the CNN into several parts. Sun et al. [9] partially addressed the challenges posed by incomplete images and spatial misalignment. They achieved this by estimating the overlapping regions between two pedestrian images through the utilization of visibility-aware part-level features. The second approach is to represent pedestrians by global features of characters. Luo et al. [10] suggested various approaches to enhance the recognition efficacy of global features. Nonetheless, these methods that rely on global features exhibit limited performance in scenarios involving occlusion and alterations in the pedestrian’s posture. The third approach is to combine local and global features, which are used to obtain a more significant representation of pedestrian features. Wang et al. [11]

proposed Multigranular Network (MGN) in order to combine fine-grained local features with global features. Park et al. [12] introduced an RRID network that not only integrates global and local features but also leverages the interplay between different body parts. Combining a pre-trained backbone with a well-designed loss function is the popular pipeline, among which triplet loss [6] and cross-entropy loss [7] are the most widely used. To enhance the integration of cross-entropy loss and triplet loss, Luo et al. [8] introduced BNNeck as a designed mechanism.

B. Occluded Person Re-Identification

The ReID methods mentioned in references [6–12] assume the visibility of the entire pedestrian body, neglecting the more complex scenario of occlusion. In reality, pedestrians are frequently obscured by objects or other individuals. To address the aforementioned issues, a solution known as occluded person re-identification has been suggested. There are two main categories of mainstream occlusion re-identification methods, which are extra information-based methods and part-to-part matching methods.

The extra information-based approach uses pose estimation, segmentation, etc., to localize the human body parts. Wang et al. [16] used the poses estimation method to learn visible local features as well as topological information during training and testing phases to achieve high detection accuracy, which leads to higher computational overhead and slower inference. Miao et al. [41] employ pose landmarks to separate valuable information from occlusion noise. Despite the progress achieved by incorporating pose landmarks, the inability to train pose-guided region extraction and the constraints imposed by predefined landmarks visibility continue to hinder matching performance.

Part-to-part matching is an approach to solving the occlusion problem by comparing the similarity between local features. In their study, Zhang et al. [13] used an approach to achieve feature alignment in their study by finding the shortest paths between local features. Sun et al. [14] proposed a visibility-aware part model(VPM), which uses self-supervised learning to learn the visibility of a perceptual region. Jia et al. [15] introduced the Mask over Similarity (MoS) method in their study, which uses Jaccard similarity to measure the similarity between character images. Specifically, they reformulated the character recognition problem due to occlusion as an unaligned set-matching problem. However, this method always suffers from misalignment, missing parts, and background confusion.

In our work, the topological prior is introduced only in the training phase to restrict a more accurate local pooling region. In addition, we work on constructing more fine-grained and robust body part features through our designed focuser and part triplet loss.

III. METHODS

The structure of our proposed part-attention based model is shown in Fig. 2. It consists of two core modules, the part

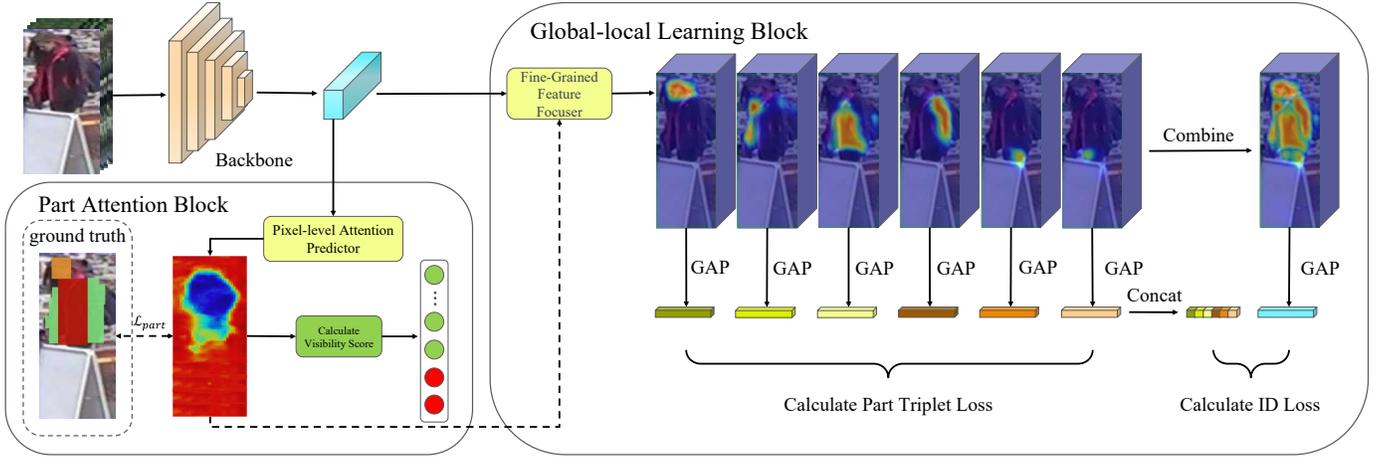


Fig. 2. The detailed structure of PAB-ReID is shown in the figure above. The PAB-ReID model comprises a part attention block responsible for generating part attention maps and a global-local learning block for extracting body part features. This paper generates six feature vectors representing the head, left hand, right hand, forehead, left leg, and right leg. In the above figure, the pedestrian’s legs are occluded, and the visibility of the occluded leg is set to 0 when calculating the visibility score.

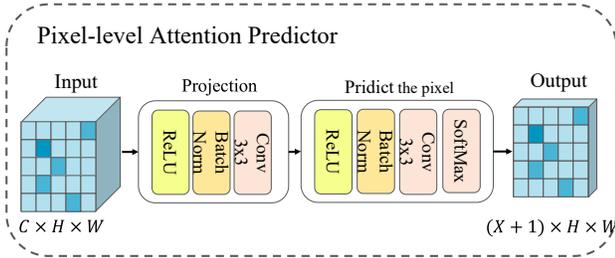


Fig. 3. The detail architecture of pixel-level attention predictor.

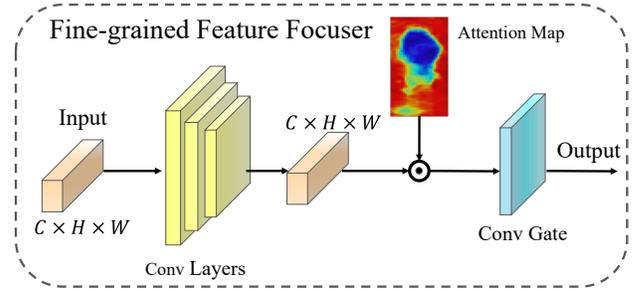


Fig. 4. The detail architecture of fine-grained feature focuser is shown in the figure above.

attention block and the global-local learning block. In the part attention block, we introduce the human parsing labels to guide the generation of the part attention maps. Meanwhile, we filter irrelevant background information, and generate fine-grained body part features in the global-local learning block.

A. Part Attention Block

As shown in Fig.2, the part attention block processes the features extracted by the backbone and generates a set of part attention maps that highlight specific body parts. In this block, we apply human parsing label to supervise the generation of attention maps for human body parts. This supervised labels consists of a number of roughly constrained regions, rather than pixel-level supervision. This is because the attention maps are also supervised by ReID-related losses, such as the ID loss and the Triplet loss, and thus the attention maps will focus on discriminative representation regions. In addition, we use part attention maps to compute the visibility of each body part. In the inference phase, we utilize only the visible body parts for the ReID task.

a) *Pixel-level Attention predictor*: The input feature map extracted by backbone is denoted as $B \in R^{C \times H \times W}$, which is subsequently processed by the pixel-level attention predictor.

The process of the predictor is shown below:

$$F_1 = \text{ReLU}(f(\delta_3(B))) \quad (1)$$

$$F = \text{Softmax}(\text{ReLU}(f(\delta_3(F_1)))) \quad (2)$$

Where δ_3 and $\text{ReLU}(\cdot)$ denotes 3×3 convolution operations and relu activation function, f is the batch normalization, the softmax function is denoted as $\text{Softmax}(\cdot)$. $F \in R^{(x+1) \times H \times W}$ denotes X part attention maps and a background attention map. In this paper, X takes the value of 6.

b) *Part Attention Loss*: We use part attention loss \mathcal{L}_{part} to supervise the pixel-level attention predictor and the loss calculation process is as follows:

$$\mathcal{L}_{part} = - \sum_{x=0}^X \sum_{w=0}^{W-1} \sum_{h=0}^{H-1} \varepsilon \cdot \log(F_x(w, h)) \quad (3)$$

$$\varepsilon = \begin{cases} \frac{\theta}{N} & \text{where } L(w, h) = x \\ 1 - \frac{N-1}{N} & \text{others} \end{cases} \quad (4)$$

Where N is the batch size and θ is the label regularisation rate. If the pixel position (w, h) of the human parsing label is attributed to the x -th body parts, the value of its pixel value $L(w, h)$ is set to x and the background position is set to 0.

c) *Visibility Score*: PAB-ReID computes a visibility score for each body part. We use 1 to denote visible parts and 0 to denote invisible parts. If at least one pixel value in a part attention map F_x , $x \in \{1, \dots, X\}$ is higher than the set value μ (μ is empirically set to 0.5), the visibility score will be set to 1.

$$V_s = \begin{cases} 1 & \text{where } \max(F_x(w, h)) > \mu \\ 0 & \text{others} \end{cases} \quad (5)$$

B. Global-local Learning Block

In the global-local learning block, our proposed fine-grained feature fuser applies part attention maps to the deep features extracted by backbone to obtain fine-grained body part features, after that gated convolution is introduced to filter the background clutter.

a) *Fine-Grained Feature Fuser*: The Pixel-level attention predictor generates X attention maps highlighting the corresponding X body parts. We first merge the X body part attention maps to generate a foreground attention map $F_f \in R^{1 \times H \times W}$:

$$F_f = \text{Concat}(F_1, \dots, F_x) \quad (6)$$

For the depth feature B extracted by backbone, we further extract the feature representation of the whole body:

$$K_1 = \text{Conv}(B) \quad (7)$$

After that, we apply the B body part attention maps as well as the foreground attention map to the overall feature representation K_1 :

$$P_i = K_i \odot F_i, \quad i \in \{f, 1, \dots, X\} \quad (8)$$

P_i , $i \in \{f, 1, \dots, X\}$ denotes the fine-grained body part features and foreground feature. Moreover, we use gated convolution to filter the background contained in the features:

$$Q_i = \text{Convgate}(F_i), \quad i \in \{f, 1, \dots, X\} \quad (9)$$

b) *global – local ReID Loss*: To increase the robustness of the model to non-discriminative features, we design the part triplet loss, and we denote the distance between two instances by the average distance of all body parts:

$$d_{parts}^{i,j} = \frac{1}{X} \sum_{x=1}^X \text{dist}_{eucl}(f_x^i, f_x^j) \quad (10)$$

Where $\text{dist}_{eucl}(f_x^i, f_x^j)$ refers to the Euclidean distances of the two parts f_x^i, f_x^j , and subsequently, the losses are calculated using the average distances of the hardest positive and hardest negative parts d_{parts}^{ap} and d_{parts}^{an} , respectively.

$$\mathcal{L}_{tri}(f_0^a, \dots, f_X^a) = [d_{parts}^{ap} - d_{parts}^{an} + \alpha] + \quad (11)$$

Where α is the part triplet loss margin, our proposed loss function makes it possible for model to focus on the most robust and discriminative parts during training, mitigating the effects of non-discriminative local and occluded features. The ablation experiments for the loss function in Fig. 5 show that part triplet loss is more effective than normal triplet loss.

In the ID classifier, our goal is to classify people with different identities, given instance x_i , with the identity label D_i . The loss function is calculated as follows:

$$\mathcal{L}_{ID} = - \sum_{i=1}^N \log \left(\frac{y(x_i, D_i)}{\sum_{j=1}^{N_{ID}} y(x_i, D_j)} \right) \quad (12)$$

N and N_{ID} are the number of instances and body points, respectively, and $y(x_i, D_i)$ denotes the probability of predicting x_i as D_i .

C. Training Procedure

The overall loss function to be optimized in the training phase is as follows:

$$\mathcal{L}_{sum} = \mathcal{L}_{tri} + \mathcal{L}_{ID} + \gamma_{part} \mathcal{L}_{part} \quad (13)$$

Where \mathcal{L}_{part} is the loss function used in the part attention block for supervised attention map generation, and γ_{part} is used to control the contribution of \mathcal{L}_{part} to the total loss, usually set to 0.35.

IV. EXPERIMENTS

A. Datasets and Settings

We conducted experiments on three mainstream occluded ReID datasets: Occluded-Duke [21], Occluded-reID [22], and P-DukeMTMC [22]. Occluded-Duke [21] is constructed based on the DukeMTMC [23] dataset. It comprises 15,618 training images for 702 identities, 2,210 occlusion query images for 519 identities, and 17,661 gallery images. There is a rich variety of variations in Occluded-Duke [21], including different viewpoints and a wide variety of obstacles, including cars, bicycles, trees, and other people. The Occluded-reID [22] dataset was a new dataset captured by a mobile camera device containing 2,000 images of 200 people. For each identity, there are five full-body images of the person and five images of the person with different types of severe occlusions. P-DukeMTMC [22] is another subset of DukeMTMC [23]. This dataset includes 12,927 training images from 665 identities, 2,163 query images from 634 identities, and 9,053 gallery images. We also evaluate our model on the DukeMTMC-ReID [23] and Market-1501 [24] datasets. The DukeMTMC-ReID [23] dataset contains 36,411 images of 1,812 pedestrians. Of these, 1,404 pedestrians were captured by more than two cameras, while 408 pedestrians were captured by only one camera. The Market-1501 [24] dataset contains 32,668 images of 1,501 individuals from 6 cameras. We use rank-k (R@K) and mean average precision (mAP) for model evaluation.

B. Implementations Details

We use Resnet-50 [20] after pre-training on ImageNet1K as the backbone network, and we remove the last down sampling of ResNet-50 in order to be able to extract deep features better. We use the Openpifpaf [42] and Mask R-CNN [43] to generate the human parsing labels. For data enhancement, images were first enhanced by random cropping and pixel filling, followed by random erasure with a probability of 0.5. Each training

TABLE I
COMPARISON OF METHODS ON OCCLUDED DATASETS

Method	Venue	Occluded Datasets					
		Occluded-Duke		Occluded-reID		P-DukeMTMC	
		Rank-1	mAP	Rank-1	mAP	Rank-1	mAP
PVPM[33]	CVPR 20	47	37.7	-	-	85.1	69.9
HG[34]	BMVC 21	61.4	50.5	-	-	-	-
PAT[35]	CVPR 21	64.5	53.6	81.6	72.1	-	-
SSGR[36]	CVPR 21	69	57.2	78.5	72.9	-	-
OAMN[25]	CVPR 21	62.6	46.1	-	-	-	-
FED[27]	CVPR 22	68.1	56.4	86.3	79.3	83.1	80.5
FRT[28]	TIP 22	70.7	61.3	80.4	71	-	-
RFCnet[29]	TPAMI 22	63.9	54.5	-	-	63.9	54.5
HCGA[32]	TIP 23	70.2	-	87.2	-	-	-
BPBreID[30]	WACV 23	66.7	54.1	76.9	68.6	91	77.8
CAAO[31]	TIP 23	68.5	59.5	87.1	83.4	92.5	81.4
QPM[37]	TMM 23	66.7	53.3	-	-	90.7	75.3
Ours		72.6	63.5	87.4	87.1	93.1	83.2

batch consisted of 32 samples of 8 IDs with 2 pictures per ID. Our model was trained in an end-to-end manner on 2 NVIDIA RTX 3090 GPUs for a total of 120 epochs, using the Adam optimizer. The learning rate increases linearly from 3.5×10^{-5} to 3.5×10^{-4} after the first ten epochs and then decays to 3.5×10^{-5} and 3.5×10^{-6} at the 40th and 70th epochs, respectively. The marginal α of the ternary loss was set to 0.3, and the label smoothing regularisation rate e was set to 0.1.

C. Comparison with the State-of-the Art Method

As shown in Table 1, PAB-ReID exhibits advanced performance on the occluded dataset, achieving 72.6%, 87.4%, and 93.1% Rank-1 and 63.5%, 87.1%, and 83.2% mAP on Occluded-Duke, Occluded-reID, and P-DukeMTMC, respectively. compared with the cutting-edge methods in terms of Rank-1 by 1.9%, 0.2%, 0.6% and mAP by 2.2%, 3.7%, 1.8% respectively. In the Table 1, PAT [9] uses a local feature generation method, and due to the lack of a priori human topology information, the model is susceptible to interference from partial omissions and misalignments, leading to less effective results than PAB-ReID. Compared to the globally-based occlusion ReID approach HG [17], PAB-ReID also demonstrates notable performance, showing an 11.2% improvement in Rank-1 accuracy and a 13% increase in mean average precision. This highlights the effectiveness of the partially-based method in addressing occlusion challenges, as the global method cannot achieve part-to-part matching. The Occluded-ReID dataset does not have a training set, our method still achieves excellent performance, which demonstrates that PAB-ReID possesses better domain adaptation.

We assess the performance of the models presented in this paper by comparing them with various cutting-edge ReID approaches, including FRT [20], RFCnet [21], HCGA [24], BPBreID [22], CAAO [23], QPM [37], etc., on the plain ReID and occluded ReID datasets, respectively, where we compare PAB-ReID with RFCnet [12], BPBreID [22], CAAO [23], and HCGA [24] models on the common ReID dataset. As shown in Table 2, our proposed method achieves state-of-

the-art performance on Market-1501 and DukeMTMC-ReID datasets. Specifically, our proposed method achieves 96.1% Rank-1 and 89.5% mAP on the Market-1501 dataset, which improves Rank-1 and mAP by 0.9% and 1.1%, respectively, compared to the existing state-of-the-art methods. Rank-1 of 91.2% and mAP of 82.5% were achieved on the DukeMTMC-ReID dataset. The above results show that PAB-ReID is also very effective when facing common person re-identification scenarios.

TABLE II
COMPARISON OF PAB-REID WITH THE STATE OF THE ART METHODS IN NORMAL REID DATASETS

Method	Normal Datasets			
	Market-1501		DukeMTMC-ReID	
	Rank-1	mAP	Rank-1	mAP
OAMN[25]	93.2	79.8	86.3	72.6
PAT[26]	95.4	88	88.8	78.2
FED[27]	95	86.3	89.4	78
FRT[28]	95.5	88.1	90.5	81.7
RFCnet[29]	95.5	89.2	90.7	80.7
BPBreID[30]	95.1	87	89.6	78.3
CAAO[31]	95.3	88	89.8	80.9
HCGA[32]	95.2	88.4	-	-
Ours	96.1	89.5	91.2	82.5

D. Ablation Studies

As shown in Table 3, we have taken ablation experiments for different modules in Occluded-Duke. The first line of Table 3 shows the performance when the part attention block is not included, the second line shows the performance when the fine-grained feature focuser is removed, the third line shows the performance when the pixel-level feature is not used predictor, the third line is the performance without using pixel-level attention predictor, and the fourth line is the ablation experiment using normal triplet loss.

a) *Part Attention Block and Pixel-level Attention predictor*: We conducted many ablation experiments for part attention block with pixel-level attention predictor to demonstrate the effectiveness of part attention block with pixel-level attention predictor. As shown in the first and third rows of

Table 3, when removing the part attention block in the PAB-ReID, the Rank-1 of the model decreases by 7.2%, and the mAP decreases by 7.3%. It proves that using human parsing label to generate body part attention maps can better guide the model to learn the features of different body parts. When we remove the pixel-level attention predictor in the part attention block, the Rank-1 and mAP of the model decrease by 1.2%

TABLE III
ABLATION EXPERIMENTS OF PAB-REID

Ablation Studies	Occluded-Duke	
	Rank-1	mAP
w/o Part attention block	65.4	56.2
w/o Fine-Grained Feature Focuser	67.3	58.3
w/o Pixel-level Attention predictor	71.4	61.7
w/o Part Triplet Loss	67.1	57.1
ALL	72.6	63.5

and 1.8%, respectively, indicating that the introduction of the pixel-level attention predictor enables the model to generate more discriminative part attention maps.

b) *Fine-Grained Feature Focuser*: We conducted experiments on the Occluded-Duke dataset to demonstrate its effectiveness for fine-grained feature focuser. As shown in the third row of Table 3, the Rank-1 and mAP of the fine-grained feature focuser model decreased by 5.3% and 5.2%, respectively, when removed from PAB-ReID. The above results show that the fine-grained feature focuser learns fine-grained the body part and foreground information by mapping the partial attention to holistic features while suppressing background information through the gated convolution.

c) *Part Triplet Loss*: Part Triplet Loss guides the model to learn fine-grained body part ReID features. We conduct comparative experiments using normal id loss, triplet loss, and id + triplet loss on the Occluded-Duke dataset to showcase the efficacy of part triplet loss. As shown in Fig. 2, when using only triplet loss on top of PAB-ReID, Rank-1, and mAP are 67.1% and 57.1%, respectively, decreased by 5.5% and 6.4% compared to using part triplet loss. The above results show that part triplet loss optimizes the intra/inter-class distance and enhances the learning of robust features.

d) *Hyper-Parameter Sensitivity Experiments*: We introduced two important hyperparameters μ and γ_{part} in PAB-ReID, the combination of which determines the model's performance. In order to explore the optimal combination of μ and γ_{part} , we conducted an ablation experiment of the hyperparameters on the Occluded-Duke dataset, and the results of the experiment are shown in Fig. 6. We first set μ and γ_{part} to 0.15 and 0.3, and then increased them by 0.2, respectively. The experimental results show that when the visibility score is larger than 0.7, the model takes part of the visible features as occluded features, which reduces the model performance and is used to regulate the proportion of \mathcal{L}_{part} in the \mathcal{L}_{sum} . When $\gamma_{part}=0.55$, the proportion of \mathcal{L}_{part} is larger, the model pays more attention to the generation of the attention graph in the part attention block and reduces

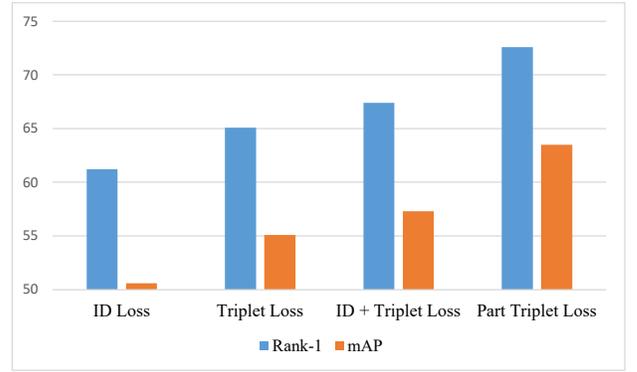


Fig. 5. Ablation experiments of loss function. The first column represents ID loss, the second column represents triplet loss, the third column represents the sum of ID and triplet loss, and the fourth column represents part triplet loss. From the above figure, it can be observed that part triplet loss can better guide the model learning.

the attention to the ReID task, resulting in a decrease in the Rank-1 of the model.

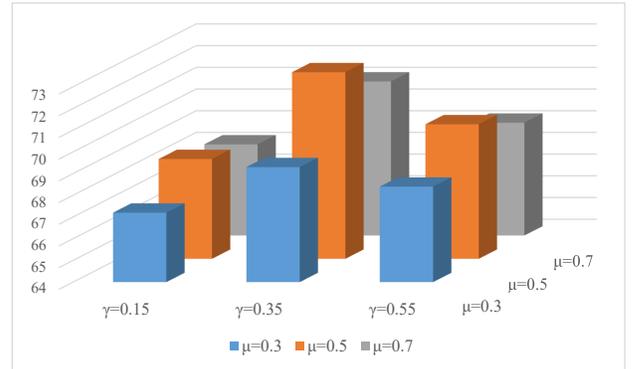


Fig. 6. The Rank-1 results from hyperparameter sensitivity experiments. The model performs best when μ and γ_{part} are set to 0.5 and 0.35, respectively.

V. CONCLUSION

In this paper, we propose a novel part-attention based model (PAB-ReID) to address the challenges in occluded person re-identification. Firstly, we design the part attention block which utilize the external human semantic information to generate ReID-related body part attention maps. This part attention maps provide more accurate feature extraction regions for human body parts. Secondly, we also propose a fine-grained feature focuser for obtaining more granular pedestrian features while suppressing the background information. Thirdly, we proposed the part triplet loss to supervise the learning of body part feature, which enhances the robustness of the model to similar body part appearance. Finally, experiments on five rigorous datasets surface our method outperforming existing state-of-the-art methods.

ACKNOWLEDGMENT

- Chen would like to thank Ge for being his guiding light in his research.

- Chen would like to thank his parents and friends for supporting his scientific endeavors.
- Supported by Promoting the Classification and Development of Colleges and Universities-Student Innovation and Entrepreneurship Training Programme Project-School of Computer (5112410852).

REFERENCES

- [1] Hou, Ruibing, et al. "Feature completion for occluded person re-identification." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44.9 (2021): 4894-4912.
- [2] Yang, Zizheng, et al. "Unleashing potential of unsupervised pre-training with intra-identity regularization for person re-identification." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.
- [3] Dai, Yongxing, et al. "Idm: An intermediate domain module for domain adaptive person re-id." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021.
- [4] Luo, Hao, et al. "Bag of tricks and a strong baseline for deep person re-identification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 2019.
- [5] He, Shuting, et al. "Transreid: Transformer-based object re-identification." *Proceedings of the IEEE/CVF international conference on computer vision*. 2021.
- [6] Liu H, Feng J, Qi M, et al. End-to-end comparative attention networks for person re-identification[J]. *IEEE transactions on image processing*, 2017, 26(7): 3492-3506.
- [7] Liu, Hao, et al. "End-to-end comparative attention networks for person re-identification." *IEEE transactions on image processing* 26.7 (2017): 3492-3506.
- [8] Luo H, Gu Y, Liao X, et al. Bag of tricks and a strong baseline for deep person re-identification[C]//*Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 2019: 0-0.
- [9] Sun Y, Xu Q, Li Y, et al. Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification[C]//*Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019: 393-402.
- [10] Luo H, Gu Y, Liao X, et al. Bag of tricks and a strong Dbaseline for deep person re-identification[C]//*Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 2019: 0-0.
- [11] Wang G, Yuan Y, Chen X, et al. Learning discriminative features with multiple granularities for person re-identification[C]//*Proceedings of the 26th ACM international conference on Multimedia*. 2018: 274-282.
- [12] Park H, Ham B. Relation network for person re-identification[C]//*Proceedings of the AAAI conference on artificial intelligence*. 2020, 34(07): 11839-11847.
- [13] Zhang, Xuan, et al. "Alignedreid: Surpassing human-level performance in person re-identification." *arXiv preprint arXiv:1711.08184* (2017).
- [14] Sun, Yifan, et al. "Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.
- [15] Jia, Mengxi, et al. "Matching on sets: Conquer occluded person re-identification without alignment." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. No. 2. 2021.
- [16] Wang, Guan'an, et al. "High-order information matters: Learning relation and topology for occluded person re-identification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [17] He, Lingxiao, and Wu Liu. "Guided saliency feature learning for person re-identification in crowded scenes." *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*. Springer International Publishing, 2020.
- [18] Song, Chunfeng, et al. "Mask-guided contrastive attention model for person re-identification." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- [19] Wang, Guan'an, et al. "High-order information matters: Learning relation and topology for occluded person re-identification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [20] He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
- [21] Miao, Jiayu, et al. "Pose-guided feature alignment for occluded person re-identification." *Proceedings of the IEEE/CVF international conference on computer vision*. 2019.
- [22] Zhuo, Jiakuan, et al. "Occluded person re-identification." *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018.
- [23] Zheng, Zhedong, Liang Zheng, and Yi Yang. "Unlabeled samples generated by gan improve the person re-identification baseline in vitro." *Proceedings of the IEEE international conference on computer vision*. 2017.
- [24] Zheng, Liang, et al. "Scalable person re-identification: A benchmark." *Proceedings of the IEEE international conference on computer vision*. 2015.
- [25] Chen, Peixian, et al. "Occlude them all: Occlusion-aware attention network for occluded person re-id." *Proceedings of the IEEE/CVF international conference on computer vision*. 2021.
- [26] Li, Yulin, et al. "Diverse part discovery: Occluded person re-identification with part-aware transformer." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021.

- [27] Wang, Zhikang, et al. "Feature erasing and diffusion network for occluded person re-identification." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022.
- [28] Xu, Boqiang, et al. "Learning feature recovery transformer for occluded person re-identification." IEEE Transactions on Image Processing 31 (2022): 4651-4662.
- [29] Hou, Ruibing, et al. "Feature completion for occluded person re-identification." IEEE Transactions on Pattern Analysis and Machine Intelligence 44.9 (2021): 4894-4912.
- [30] Somers, Vladimir, Christophe De Vleeschouwer, and Alexandre Alahi. "Body part-based representation learning for occluded person Re-Identification." Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. 2023.
- [31] Zhao, Cairong, et al. "Content-Adaptive Auto-Occlusion Network for Occluded Person Re-Identification." IEEE Transactions on Image Processing (2023).
- [32] Dou, Shuguang, et al. "Human Co-Parsing Guided Alignment for Occluded Person Re-Identification." IEEE Transactions on Image Processing 32 (2022): 458-470.
- [33] Gao, Shang, et al. "Pose-guided visible part matching for occluded person reid." Proceedings of the IEEE/CVF conference
- [34] Kiran, Madhu, et al. "Holistic guidance for occluded person re-identification." arXiv preprint arXiv:2104.06524 (2021).
- [35] Li, Yulin, et al. "Diverse part discovery: Occluded person re-identification with part-aware transformer." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.
- [36] Yan, Cheng, et al. "Occluded person re-identification with single-scale global representations." Proceedings of the IEEE/CVF international conference on computer vision. 2021.
- [37] Wang, Pengfei, et al. "Quality-aware part models for occluded person re-identification." IEEE Transactions on Multimedia (2022).
- [38] Sun, Yifan, et al. "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)." Proceedings of the European conference on computer vision (ECCV). 2018.
- [39] Zhang, Zhong, Haijia Zhang, and Shuang Liu. "Person re-identification using heterogeneous local graph attention networks." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2021.
- [40] Zhu, Kuan, et al. "Identity-guided human semantic parsing for person re-identification." Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16. Springer International Publishing, 2020.
- [41] Miao, Jiayu, et al. "Pose-guided feature alignment for occluded person re-identification." Proceedings of the IEEE/CVF international conference on computer vision. 2019.
- [42] Kreiss, Sven, Lorenzo Bertoni, and Alexandre Alahi. "Openpipaf: Composite fields for semantic keypoint detection and spatio-temporal association." IEEE Transactions on Intelligent Transportation Systems 23.8 (2021): 13498-13511.
- [43] He, Kaiming, et al. "Mask r-cnn." Proceedings of the IEEE international conference on computer vision. 2017.
- [44] Peng, Yunjie, et al. "Deep Learning Based Occluded Person Re-Identification: A Survey." ACM Transactions on Multimedia Computing, Communications and Applications 20.3 (2023): 1-27.
- [45] Jia, Mengxi, et al. "Semi-attention partition for occluded person Re-identification." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 37. No. 1. 2023.
- [46] Dong, Neng, et al. "Erasing, transforming, and noising defense network for occluded person re-identification." IEEE Transactions on Circuits and Systems for Video Technology (2023).
- [47] He, Shuting, et al. "Region generation and assessment network for occluded person re-identification." IEEE Transactions on Information Forensics and Security (2023).
- [48] Liu, Wenfeng, et al. "Learning Occlusion Disentanglement with Fine-grained Localization for Occluded Person Re-identification." Proceedings of the 31st ACM International Conference on Multimedia. 2023.
- [49] Gao, Zan, et al. "A Semantic Perception and CNN-Transformer Hybrid Network for Occluded Person Re-identification." IEEE Transactions on Circuits and Systems for Video Technology (2023).
- [50] Zhao, Cairong, et al. "Content-Adaptive Auto-Occlusion Network for Occluded Person Re-Identification." IEEE Transactions on Image Processing (2023).
- [51] Zhao, Cairong, et al. "Content-Adaptive Auto-Occlusion Network for Occluded Person Re-Identification." IEEE Transactions on Image Processing (2023).
- [52] Yan, Gang, et al. "Part-based Representation Enhancement for Occluded Person Re-identification." IEEE Transactions on Circuits and Systems for Video Technology (2023).
- [53] Huang, Meiyan, et al. "Reasoning and Tuning: Graph Attention Network for Occluded Person Re-Identification." IEEE Transactions on Image Processing 32 (2023): 1568-1582.
- [54] Zhou, Shuren, and Mengsi Zhang. "Occluded person re-identification based on embedded graph matching network for contrastive feature relation." Pattern Analysis and Applications 26.2 (2023): 487-503.
- [55] R. Guan, Z. Li, W. Tu, J. Wang, Y. Liu, X. Li, C. Tang, and R. Feng, "Contrastive multi-view subspace clustering of hyperspectral images based on graph convolutional networks," IEEE Transactions on Geo-science and Remote Sensing, vol. 62, pp. 1–14, 2024.
- [56] R. Guan, Z. Li, X. Li, and C. Tang, "Pixel-superpixel contrastive learning and pseudo-label correction for hyperspectral image clustering," arXiv preprint

- arXiv:2312.09630, 2023
- [57] R. Guan, Z. Li, T. Li, X. Li, J. Yang, and W. Chen, "Classification of heterogeneous mining areas based on rescapsnet and gaofen-5 imagery," *Remote Sensing*, vol. 14, no. 13, p. 3216, 2022.
- [58] J. Liu, R. Guan, Z. Li, J. Zhang, Y. Hu, and X. Wang, "Adaptive multi-feature fusion graph convolutional network for hyperspectral image classification," *Remote Sensing*, vol. 15, no. 23, p. 5483, 2023.
- [59] W. Tu, R. Guan, S. Zhou, C. Ma, X. Peng, Z. Cai, Z. Liu, J. Cheng, and X. Liu, "Attribute-missing graph clustering network," in *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI)*, 2024.
- [60] Lai, Songning, et al. "Shared and private information learning in multimodal sentiment analysis with deep modal alignment and self-supervised multi-task learning." arXiv preprint arXiv:2305.08473 (2023).
- [61] Lai, Songning, et al. "Multimodal sentiment analysis: A survey." *Displays* (2023): 102563.
- [62] Lai, Songning, et al. "Predicting Lysine Phosphoglyceration Sites using Bidirectional Encoder Representations with Transformers & Protein Feature Extraction and Selection." *2022 15th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*. IEEE, 2022.
- [63] Lai, Songning, et al. "Faithful Vision-Language Interpretation via Concept Bottleneck Models." *The Twelfth International Conference on Learning Representations*. 2023.
- [64] Xu, Haoxuan, et al. "Cross-domain car detection model with integrated convolutional block attention mechanism." *Image and Vision Computing* 140 (2023): 104834.
- [65] Zhan, H., Zhang, K., Hu, C., & Sheng, V. S. (2022). New threats to privacy-preserving text representations. In *55th Annual Hawaii International Conference on System Sciences, HICSS 2022* (pp. 768-777). IEEE Computer Society.
- [66] Zhan, H., Zhang, K., Lu, K., & Sheng, V. S. (2023, September). Measuring the privacy leakage via graph reconstruction attacks on simplicial neural networks (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 37, No. 13, pp. 16380-16381).
- [67] Zhan, H., Zhang, K., Chen, Z., & Sheng, V. S. (2023, October). Simplex2vec Backward: From Vectors Back to Simplicial Complex. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management* (pp. 4405-4409).
- [68] Zhan, H., Gao, L., Zhang, K., Chen, Z., & Sheng, V. S. (2023, October). Defending the Graph Reconstruction Attacks for Simplicial Neural Networks. In *2023 IEEE 10th International Conference on Data Science and Advanced Analytics (DSAA)*(pp. 1-9). IEEE.
- [69] Zhan, H., Zhang, K., Hu, C., & Sheng, V. (2021, October). Multi-objective privacy-preserving text representation learning. In *Proceedings of the 30th acm international conference on information & knowledge management* (pp. 3612-3616).
- [70] Chen, Zhihao, and Yiyuan Ge. "Occluded Cloth-Changing Person Re-Identification." arXiv preprint arXiv:2403.08557 (2024).
- [71] Ge, Yiyuan, et al. "Lightweight Traffic Sign Recognition Model Based on Dynamic Feature Extraction." *International Conference on Applied Intelligence*. Singapore: Springer Nature Singapore, 2023.
- [72] Ge, Yiyuan, et al. "End-to-End Person Search Based on Content Awareness." *2023 International Conference on Image Processing, Computer Vision and Machine Learning (ICICML)*. IEEE, 2023.
- [73] Zhang, Ji, et al. "An Efficient Convolutional Multi-Scale Vision Transformer for Image Classification." *2023 International Conference on Image Processing, Computer Vision and Machine Learning (ICICML)*. IEEE, 2023.
- [74] Chen, Zhihao, et al. "Multi-branch Person Re-identification Net." *2023 International Conference on Image Processing, Computer Vision and Machine Learning (ICICML)*. IEEE, 2023.
- [75] Tian, Ye, et al. "View while Moving: Efficient Video Recognition in Long-untrimmed Videos." *Proceedings of the 31st ACM International Conference on Multimedia*. 2023.
- [76] Mengyu Yang, et al. "AdaViPro: Region-based Adaptive Visual Prompt for Large-Scale Models Adapting." arXiv preprint arXiv:2403.13282 (2024).
- [77] Chen Z, Ge Y. MambaUIE&SR: Unraveling the Ocean's Secrets with Only 2.8 FLOPs[J]. arXiv preprint arXiv:2404.13884, 2024.
- [78] Chen Z, Ge Y. Part-Attention Based Model Make Occluded Person Re-Identification Stronger[J]. arXiv preprint arXiv:2404.03443, 2024.

This figure "fig1.png" is available in "png" format from:

<http://arxiv.org/ps/2404.03443v4>