# Region Generation and Assessment Network for Occluded Person Re-Identification

Shuting He, Weihua Chen, Kai Wang, Hao Luo, Fan Wang, Wei Jiang, Henghui Ding

*Abstract*—**Person Re-identification (ReID) plays a more and more crucial role in recent years with a wide range of applications. Existing ReID methods are suffering from the challenges of misalignment and occlusions, which degrade the performance dramatically. Most methods tackle such challenges by utilizing external tools to locate body parts or exploiting matching strategies. Nevertheless, the inevitable domain gap between the datasets utilized for external tools and the ReID datasets and the complicated matching process make these methods unreliable and sensitive to noises. In this paper, we propose a Region Generation and Assessment Network (RGANet) to effectively and efficiently detect the human body regions and highlight the important regions. In the proposed RGANet, we first devise a Region Generation Module (RGM) which utilizes the pre-trained CLIP to locate the human body regions using semantic prototypes extracted from text descriptions. Learnable prompt is designed to eliminate domain gap between CLIP datasets and ReID datasets. Then, to measure the importance of each generated region, we introduce a Region Assessment Module (RAM) that assigns confidence scores to different regions and reduces the negative impact of the occlusion regions by lower scores. The RAM consists of a discrimination-aware indicator and an invariance-aware indicator, where the former indicates the capability to distinguish from different identities and the latter represents consistency among the images of the same class of human body regions. Extensive experimental results for six widely-used benchmarks including three tasks (occluded, partial, and holistic) demonstrate the superiority of RGANet against state-of-the-art methods.**

*Index Terms*—**Person Re-Identification, Region Generation and Assessment Network (RGANet), Region Generation Module (RGM), Region Assessment Module (RAM).**
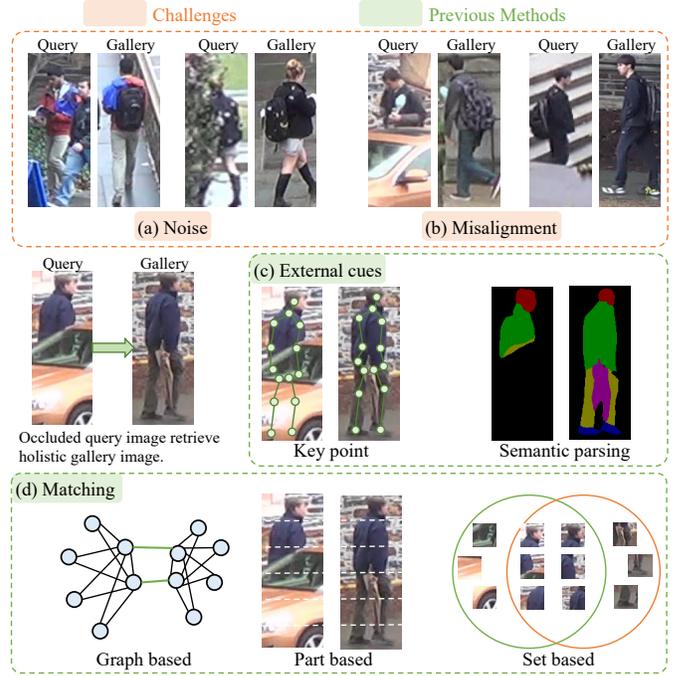


Fig. 1. The upper block illustrates the two key challenges in occluded ReID. (a) Noises caused by occlusions interfere the feature extraction. (b) The occluded image containing only part of the human body results in spatial misalignment. The lower block shows two mainstreams to solve the challenges. (c) External cues are utilized to detect and align body parts. (d) Matching based methods establish alignment relationships between feature patches obtained by specific rules.

## I. INTRODUCTION

**P**ERSON re-identification (ReID) aims to match images of the target person across different times, locations, and camera views. Holistic person ReID community has witnessed significant improvement by many advanced deep learning algorithms [1]–[6] and well-annotated large-scale datasets [7]–[13]. However, person ReID is still challenging due to the difficulty in extracting robust feature representations under the occlusion scenarios. As humans are occluded by clutters and obstacles commonly in practical surveillance systems, occluded person ReID [10], [14] deserves further study owing to its great value in real-world scenarios.

Compared with holistic person ReID, there exist two critical issues for occluded person ReID: (1) Occlusions (*e.g.*, cars, trees, boards, persons) bring various noises which interference the extraction of discriminative feature, as shown in Fig. 1(a). (2) The occluded image contains only part of the human body, leading to the issue of spatial misalignment, as shown in Fig. 1(b). Most existing works tackle such challenges

by locating the non-occluded body parts and aligning these body parts. These approaches can be categorized into two mainstreams, *i.e.*, external cues based methods and matching based methods.

External cues based methods utilize some off-the-shelf models, *e.g.*, semantic parsing [15]–[20] or key point models [10], [17], [21]–[24], to facilitate detecting and aligning the body part, as shown in Fig. 1(c). They generally use external cues to divide the human body into several informative parts, *e.g.* "head", "upper body", "lower body", and "foot" from top to bottom. However, the inevitable domain gap between the images utilized for training external tools and the ReID datasets significantly poses the obstacle of obtaining an effective and efficient model. Matching based strategies [1], [2], [22], [25]–[28] aim to construct the alignment relationship between local features patches obtained by specific rules as illustrated in Fig. 1(d). Nevertheless, complicated matching algorithms lead to high computational costs during the inference

stage. Especially, the computational complexity is overwhelming when handling large-scale datasets and matching methods may collapse in severely occluded scenarios.

In this work, we propose an efficient and effective **R**egion **G**eneration and **A**ssessment **Net**work (**RGANet**) that enables the network to adaptively locate human body parts and introduce confidence score to guide the network in training and matching. We introduce text clues into human part detection and propose a Region Generation Module (RGM) to produce the human body regions using semantic prototypes by textual descriptions. To this end, we utilize the large-scale pre-trained vision-language models, *e.g.*, CLIP [29] and ALIGN [30], to project semantic words (names of human body regions) to vision-aligned textual prototypes. The large-scale pre-trained vision-language models have strong capability of matching visual and textual features with great generalization ability towards open-world vocabulary. Therefore, each vision-aligned textual prototype represents the corresponding class and is discriminative enough to distinguish each other. In such a way, we convert the problem of partitioning each image properly into classifying each pixel of the image to a specific category from "head", "upper body", "lower body", and "foot". CLIP, being an external model, also exists domain gap with ReID datasets. To address this domain gap, we employ prompt learning to introduce additional learnable context tokens that are trained on the ReID datasets. This process enables it to adapt and align with the specific characteristics of the ReID task. Therefore, the domain gap between the dataset utilized for CLIP and the ReID datasets is effectively eliminated. We feed these category names into our designed prompt and produce the prototypes by CLIP's pre-trained text encoder. Then, each pixel of the ReID image is labeled by the class-specific prototypes nearest to its feature representation, by which the image feature map is segmented into different regions. Such regions are adaptively generated under the guidance of semantic clues and training objectives, and thus can better capture informative regions than fixed external cues based methods. Moreover, once the training is completed, the textual prototypes optimized after the training stage can be directly used for testing, rather than having to extract them again like external tools methods. By generating these human body regions, we can not only focus more on informative locations that are useful to differentiate different person identities apart, but also suppress background and occluded parts that are detrimental to the ReID performance.

Then, we propose a Region Assessment Module (RAM) to highlight the more importance regions. The clues contained in these detected regions show different importance in ReID and some occlusion part may have negative impact to the performance. For example, when a person's feet are occluded, the detected foot region cannot provide useful information for re-identifying this person. Besides, sometimes it is unavoidable to include some noise of the occlusion or background in the generated informative regions, which will interfere with the learning of the network. In the field of object detection, the credibility of each bounding box will be determined by a confidence score. Inspired by it, the Region Assessment Module (RAM) is introduced to evaluate the quality of the

generated regions and assign different confidence scores to these regions, which greatly helps to reduce the negative impact of the occlusion part. The proposed RAM consists of a discrimination-aware indicator and an invariance-aware indicator, which are used to measure the discrimination and invariance of the features of our generated regions, respectively. The discrimination-aware indicator uses a self-learning manner to measure the discrimination capability, *i.e.*, whether this region can help to distinguish different identities. The invariance-aware indicator aims to learn the invariant information embedded in all images from the same category of the human body regions, *i.e.*, what is consistent among the images of the same category of the human body regions. Combining these two complementary scores, the confidence score is obtained which is employed in the loss function for robust ReID feature learning and matching. With the help of RAM, the informative region features are more discriminative and invariant, which significantly enhances the feature representation capability of the regions. Finally, the final matching process can be performed by the summation of human body regions distance aggregation weighted by the confidence score.

The main contributions of this paper are summarized as follows:

- We propose a **R**egion **G**eneration and **A**ssessment **Net**work (**RGANet**) to adaptively locate human body parts using semantic information and select the more informative parts to enhance occluded person ReID.

- We propose a Region Generation Module (RGM) that generates human part regions using semantic representations of CLIP. Regions are thus endowed with a better discriminative ability to mitigate occlusion interference.

- To further analyze the generated regions, a Region Assessment Module (RAM) is introduced to measure the importance of these regions, which greatly lessens the adverse impact of occlusion or uninformative regions on the network.

- The proposed RGANet achieves state-of-the-art performance on occluded, partial, and holistic ReID benchmarks including Occluded-DukeMTMC [10], Occluded-REID [14], Partial-REID [12], Partial-iLDIS [11], Market-1501 [7], and MSMT17 [31].

## II. RELATED WORK

### A. Occluded Person Re-Identification

Compared with holistic person ReID, occluded person ReID is more challenging due to noise interference and spatial misalignment. Basically, existing approaches can be divided into two streams, external cues based methods and matching based approaches.

Previous methods utilize external models such as semantic parsing [15]–[17], [20], [32], [33], human pose [10], [17], [21]–[24] to locate aligned body parts. With the help of extra semantic information, such methods are capable of aligning parts precisely and extracting more robust feature representation. For example, To facilitate probe and gallery feature generation and matching, Miao *et al.* [10] propose a

pose-guided feature alignment method (PGFA) making use of the external human pose models. Gao *et al.* [21] design a pose-guided visible part matching algorithm (PVPM). Under the guidance of pose estimation and graph matching, it jointly extracts features and mines the visibility of parts through attention heatmaps. HOReID is devised by Wang *et al.* [22] to extract discriminative features and achieve robust alignment via building high-order relation and topology information according to the estimation of key-points. He *et al.* [17] introduce a new matching approach GASM to combine the pose and semantic mask information to output saliency heatmap which is utilized to further supervise discriminative feature learning and conduct adaptive spatial matching. Despite external cues contributing to alleviating the challenges of occluded ReID, the inevitable domain gap between the images used for training external models and the ReID datasets lead to unreliable external cues and hinder the subsequent ReID process.

Matching based strategies [1], [2], [22], [25], [26], [34], [35] aim to construct the alignment relationship between local features patches obtained by specific rules and it can be further divided into part-to-part matching [1], [2], [26], [34], [35], graph base matching [22], and set based matching [25]. Sun *et al.* [1] divide feature maps into several horizontal stripe embeddings and train non-shared classifiers on each one. Using K-means clustering algorithms, Zhu *et al.* [26] detect human body parts and additional personal belongings from the pixel level. Yao *et al.* [35] design the part loss to locate human body parts which enforce the network to extract representations for various parts. HOReID [22] and PVPM [21] both apply graph matching to construct alignment relationship between parts. Jia *et al.* [25] extract features along the channel dimension by a pattern set capturing one particular visual pattern and utilize Jaccard distance instead of Euclidean distance to compute the similarity between pattern sets. These methods extract and align local features through a self-supervision manner without external cues. Although complicated matching algorithms lead to high computational costs during the inference stage. Besides, the predicted result largely depends on the quality of obtained regions which is susceptible to noise. Differing from the existing works, our RGANet aims to explicitly extract the representation of the target person with the help of pre-trained CLIP and the learnable prompt is designed to eliminate domain gap between CLIP datasets and ReID datasets. Moreover, during inference, we no longer require external tools since we have optimized prototypes for segmenting the feature map. Besides, it is capable of awaring of the confidence of features, which lessens the interference of irrelevant parts.

### B. Visual-Language Pretraining

Nowadays, Visual-Language pretraining has gained increasing attention and reached superior performance on a large number of multi-modal downstream tasks. Several methods [29], [36], [37] utilize large-scale data sources and exploit semantic supervision to extract visual representations with text representations. MIL-NCE [36] has the capability of solving misalignments problem with narrated videos and obtaining strong video feature representations from noisy large-scale dataset HowTo100M [38]. SimVLM [37] lessens the training complexity through taking advantage of large-scale weak supervision which is trained under a single prefix language modeling function in an end-to-end manner. Recently, Contrastive Language-Image Pretraining (CLIP) [29] has obtained remarkable results in multi-modal zero-shot learning, which indicates cross-modal feature representations can be aligned in a shared embedding space. The remarkable generalization capability of CLIP is owing to large-scale training samples of about 400 million image-text pairs obtained from the Internet.

Recently, some researchers introduces the idea of prompt learning, a prevalent trend in NLP, to the vision area. Zhou *et al.* [39] introduce Context Optimization (CoOp) which obtain significant improvements over hand-crafted prompts by simply adding a set of learnable vectors in the prompt. To improve generalization ability, Zhou *et al.* [40] propose Conditional Context Optimization (CoCoOp), which extends CoOp through further appending a lightweight neural network aiming to generate an input-conditional token for each image. Our work is built upon the CLIP and leverages its multi-modal alignment capability for retrieving specific semantic parts.

### III. APPROACH

In this section, we describe the proposed framework of RGANet, which is shown in Fig. 2. The details of the framework, its two major modules, and its training and testing process will be introduced as below.

### A. Overview of RGANet

As shown in Fig. 2, we first extract the image feature $F$ by an image encoder. A global feature $F_g$ is obtained by a Global Average Pooling (GAP) on $F$. Region Generation Module (RGM) is applied over the feature map $F$ to locate the position of different classes of human body regions. Specifically, with the help of CLIP, we first obtain a text-generated prototype, and class segmentation masks of human parts is generated by labeling each pixel of $F$ as the class of the nearest prototype. Then the masked average pooling is utilized to obtain the class-specific region feature $F_j^r (j = 1, \cdots, N)$. To evaluate the confidence of the generated regions, region features are fed into Region Assessment Module (RAM). RAM consists of two parts, (1) a self-attention mechanism as discrimination-aware indicator to assign each region with a discrimination score using a fully connected layer followed by a sigmoid function; (2) an invariance-aware indicator utilizing a memory bank to compare the given region with its aligned regions obtains an invariance score.

### B. Revisiting CLIP

Contrastive Language-Image Pre-training (CLIP) [29] first extracts image and text features by an image encoder and a text encoder, respectively. Then the image features are forced to match with their corresponding text features under a contrastive loss, which enables image features and their paired text features to be aligned accordingly. After pre-training, given a
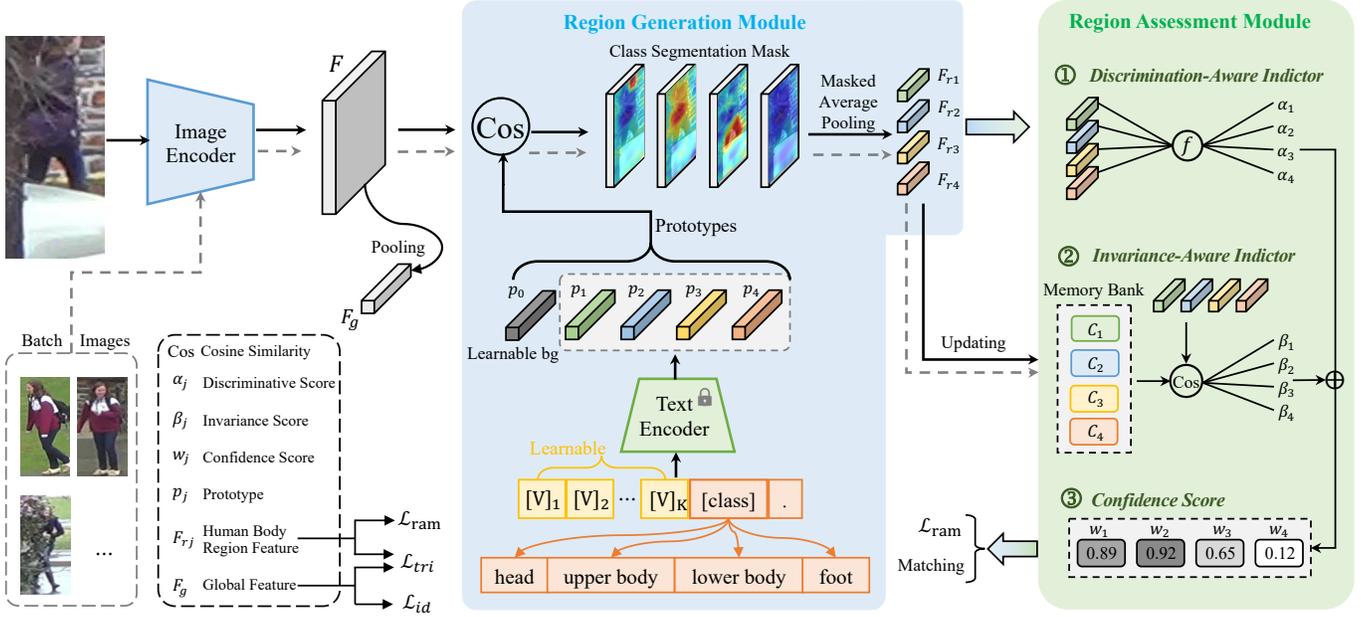
Fig. 2. The pipeline of our proposed **R**egion **G**eneration and **A**ssessment **Net**work (**RGANet**). It consists of two components: Region Generation Module (middle) and Region Assessment Module (right). The RGM employs CLIP to generate text embeddings as prototypes and then obtain class segmentation mask indicating informative regions by calculating cosine similarity between image features and these textual prototypes. RAM is comprised of the discrimination-aware indicator and the invariance-aware indicator. Discrimination-aware indicator uses $f$, self-attention layer, to measure the discrimination of features for each region. Invariance-aware indicator compares a given human body region feature $F_j^r$ with the other same category of the human body regions stored in memory bank $\mathcal{M}$. The gray dotted line indicates that the images are used to update the memory bank during the training process. Combining these two complementary scores ($\alpha_j$ and $\beta_j$), the confidence score $w_j$ is obtained for further network training and feature matching.

category name and image of an unseen category that has never been seen in the training of CLIP, CLIP can well match the textual feature obtained according to the category name and visual feature of the image, showing a great generalization ability to unseen categories.

CLIP shows strong performance in many vision tasks, for example, zero-shot learning [41]–[43] and referring image segmentation [44], [45]. A direct way to apply CLIP for zero-shot learning is adopting text embeddings generated by the text encoder of CLIP as weights of a prototype to label visual features into a specific class, i.e., $< T_i, I_i >$, where $I_i$ is image feature encoded by CLIP's image encoder, $T_i$ is prototype weights from CLIP's text encoder, and $< \cdot, \cdot >$ represents calculating cosine similarity. Inspired by it, in this work, we take advantage of the strong generalization ability of CLIP and input category names of human body regions into a text encoder of CLIP, generating semantic embeddings of the human body for classifying informative regions.

### C. Region Generation Module

To locate the informative regions and suppress the influence of occlusions/background, we introduce a Region Generation Module (RGM) that aims to adaptively capture more discriminative and well-aligned regions. The architecture of our proposed Region Generation Module is shown in the middle part of Fig. 2. In the proposed RGM, we first employ CLIP to generate text embeddings as prototypes and then obtain class segmentation masks indicating informative regions by calculating cosine similarity between image features and these textual prototypes.

**Prototype Extraction**. We assume the class names set has $N$ categories including "head", "upper body", "lower body", "foot". The class name from the class name set is placed into a learnable prompt ( introduced in **Prompt Design**) and then fed into the CLIP text encoder. Then we can obtain $N$ text embeddings, denoted as $\mathcal{P} = \{p_j \in \mathbb{R}^d | j = 1, \cdots, N\}$. In our ReID task, we also need a "background" category that indicates pixels do not belong to any human part classes. For the "background" category, one class name is insufficient to describe it. As such, we add an extra learnable embedding $p_0 \in \mathbb{R}^d$ for "background" class to represent the comprehensive analysis. Since we have obtained optimized prototypes for segmenting the feature map, we do not need to extract prototypes again during inference which is superior to the external tools based approaches.

**Class Segmentation Mask**. Meantime, we get the image feature map by our image encoder, denoted as $F \in \mathbb{R}^{d \times H \times W}$, where $d$, $H$, $W$ are channel numbers, height, and width of the feature map, respectively. $F$ has the same channel number as these textual embeddings. Each textual embedding contains category-related discriminative clues and is employed to measure how similar the pixel features in $F$ are to the text embedding. To get class segmentation masks $\mathcal{S}_j \in \mathbb{R}^{H \times W} (j = 1, \cdots, N)$ for each of the human body regions categories, we compute the cosine similarity between image feature $F$ and each of the textual embeddings and apply a softmax over the similarity maps to normalize them, i.e.,

$$\mathcal{S}_j^{(x,y)} = \frac{\exp(\gamma < F^{(x,y)}, p_j >)}{\sum_{p_k \in \{\mathcal{P} \cup p_0\}} \exp(\gamma < F^{(x,y)}, p_k >)}, \quad (1)$$

where $(x, y)$ denotes spatial position of pixels, $\gamma$ is an amplification factor, $< a, b >$ indicates calculating cosine similarity between $a$ and $b$. Each position in $S_j$ denotes the similarity between this position's visual feature with $j$-th textual embedding $p_j$. "background" is used to suppress the similarity on pixels that belong to none of the human body regions. Then, we conduct masked average pooling over image feature $F$ to obtain the global representation of each human body region feature $F_j^r$,

$$F_j^r = \frac{\sum_{x,y} S_j^{(x,y)} F^{(x,y)}}{\sum_{x,y} S_j^{(x,y)}}, \tag{2}$$

where $F_j^r \in \mathbb{R}^d (j = 1, \cdots, N)$. The proposed Region Generation Module (RGM) detects informative human body regions, by measuring feature similarity between visual features and textual features of human body regions classes provided by a pre-trained CLIP model. The generated regions emphasize body part-related positions while suppressing unrelated ones, as shown in Fig. 2. With the generated human body-related regions, we can focus more on extracting identification-related clues from these semantic regions while avoiding being affected by noisy information.

**Prompt Design**. The vanilla CLIP is not designed for person ReID and we explore designing feasible text prompts for person ReID. A direct way to use CLIP is by utilizing the hand-crafted prompts, while prompts given by CLIP are specifically suitable for image classification and may have an adverse effect on our task. Consequently, we experiment with different prompts designs to find the most useful one for our RGM, which is "a [CLASS] part of a person". Recently, learning-based prompt [39], [40] shows great promise for adapting the CLIP by prompt design on different downstream tasks. Inspired by it, we introduce learnable vectors to model prompt context for occluded ReID task. Specifically, a generalized prompt can be expressed as $[V]_1 [V]_2 \ldots [V]_K [CLASS]$ where each $[V]_k$ ($k \in \{1, \ldots, K\}$) is a vector set as learnable parameters which have the same dimensions with word embeddings (*i.e.*, 512 for CLIP), and $K$ is the number of learnable context vectors. An additional training process is indispensable to optimize these context parameters.

To supervise the prompt learning, we perform the segmentation loss $\mathcal{L}_{\text{seg}}$ as follows:

$$\mathcal{L}_{\text{seg}} = \sum_{x,y} \sum_{j=1}^{N} \mathbb{1}[\hat{\mathcal{S}}^{(x,y)} = j] \log \mathcal{S}_j^{(x,y)}, \tag{3}$$

where $\hat{\mathcal{S}}$ is the pseudo ground truth segmentation mask generated from the previous work [46]. $\mathbb{1}[\cdot]$ is an indicator math function, if the argument is true it output 1, or 0 otherwise. The incorporation of this learnable prompt effectively adapts CLIP to the ReID context. It is worth noting that hand-crafted prompts do not require additional training procedures or external pseudo-segmentation masks and still achieve decent results (please kindly refer to ablation stduy in TABLE VI for details).

### D. Region Assessment Module

The clues contained in these body part regions play different roles in ReID. Besides, when facing occlusion, the corresponding regions contain less useful information and may bring noise. If these occluded regions are fed into the training of the network as the informative regions, it will interfere with learning discriminative features, because the network may regard the same occasion as the signal of the same person wrongly. In the field of object detection, a confidence score will be assigned to each bounding box to estimate its credibility. Hence, inspired by it, we propose a Region Assessment Module (RAM) to give them different confidence scores so that they make different contributions to the network training and matching. The confidence score assigned to each region here is designed by taking these two aspects into consideration, *i.e.*, discrimination among different classes, and invariance within the same classes[1].

**Discrimination-Aware Indicator.** This indicator is designed to capture the region's contributions to each sample. It is expected that the discrimination-aware indicator assigns higher scores to the more discriminative regions and lower scores to the background or occluded regions.

The discrimination-aware indicator takes $F_j^r$ as input and outputs a discrimination score for each region. Specifically, it is composed of a fully connected (FC) layer and followed by a sigmoid activation function, which can be expressed as,

$$\alpha_j = \sigma(W_a^\top F_j^r), \tag{4}$$

where $\alpha_j$ is the discrimination score of the $j$-th region, the parameters of the FC layer are defined by $W_a \in \mathbb{R}^{d \times 1}$ and the sigmoid function is defined by $\sigma$. The discrimination-aware indicator pushes the features to learn the inherent relations of images from different identities and measures the discrimination ability of each region. It is worth noting that the parameters of the indicator are shared among different regions.

**Invariance-Aware Indicator.** The invariance-aware indicator makes sure that the features of samples within same category of the human body regions would explore the invariant characteristic shared within these regions.

We introduce a memory bank $\mathcal{M}$ to collect invariance characteristics, inspired by [47]. The key process of $\mathcal{M}$ generation includes memory initialization and memory update. Unlike [47], we initialize the memory with the class centers instead of ID centers in the training set. We utilize average features of the same category of the human body regions to act as class centers. In the early stage of network training, the features from RGM may not be extracted correctly, so we use the fixed-stripe feature to initialize the memory. It is worth noting that a forward step is executed at the beginning of the training process to initialize the class centers, which are then constantly updated as training progresses. The mean of the extracted features with the same category of the human body regions $j$ are utilized to update the $j$-th center $C_j$ in the mini-batch

$$C_j = m_u C_j + (1 - m_u) \frac{1}{|B_j|} \sum_{F_j^r \in B_j} F_j^r, \tag{5}$$

---

[1]Here classes refer to the human body region classes.

where $B_j$ represents the human body region features involved with class $j$ in the mini-batch, $m_u$ denote the momentum updating rate, $F_j^T$ is the region feature produced by RGM. And we select regions with $\alpha_j$ greater than 0.85 to update the memory bank which ensures that the features in the memory bank are high quality as possible.

Given a region tensor $F_j^r$, we take out the aligned class center $C_j$ and compute the cosine similarity between $F_j^r$ and $C_j$ to get invariance score $\beta_j$ as:

$$\beta_j = \frac{\exp(< F_j^r, C_j >)}{\sum_{j=1}^N \exp(< F_j^r, C_j >)}, \quad (6)$$

where $\beta_j$ denotes the relevance between the given region feature $F_j^r$ and $j$-th class center $C_j$. $< a, b >$ indicates calculating cosine similarity between $a$ and $b$. With the help of an invariance-aware indicator, when there is a region that is full of occlusions, it will give a small score to suppress its impact on the network after comparing it with a common well-generated region in the memory bank.

**Combination of the Indicators.** Finally, the $\alpha_j$ obtained by discrimination-aware indicator and the $\beta_j$ collected by invariance-aware indicator are summed together to get the final comprehensive confidence score $w_j$. A softmax function is then used to normalize the $w_j$ to obtain the final confidence score. We apply this confidence score $w_j$ to the loss function and in the later feature matching in Eq. (10). With the confidence score $w_j$, we multiply it by the cross-entropy loss of ReID. The loss function is described as:

$$\mathcal{L}_{\mathrm{ram}} = -w_j log(softmax(W_b^\top F_j^r)), \quad (7)$$

where $W_b$ is the classifier for $j$-th region. The $\mathcal{L}_{\mathrm{ram}}$ has a positive correlation with the $w_j$. When the quality of the given feature is not satisfactory, the network will produce a small $w_j$. Therefore, the impact on network training will be small, and vice versa. In summary, RAM combines two complementary aspects: discrimination and invariance to perceive the contribution of the generated regions and provides a comprehensive awareness of each given region.

### E. Training and Testing Process

In the training phase, we calculate cross-entropy and triplet losses [4] for both region's features and global features. The overall objective function is formulated as:

$$\mathcal{L}_{total} = \sum_j \mathcal{L}_{\mathrm{ram}}(F_j^r) + \sum_j \mathcal{L}_{tri}(F_j^r) + \mathcal{L}_{id}(F_g) + \mathcal{L}_{tri}(F_g), \quad (8)$$

where $\mathcal{L}_{id}$, $\mathcal{L}_{tri}$ denotes the cross-entropy loss, triplet losses, respectively.

In the testing phase, inspired by the [10], the distance $d_j$ of the $j$-th region between query and gallery images is:

$$d_j = D(F_j^{r,q}, F_j^{r,g})(j = 1, \ldots, N), \quad (9)$$

where $D(\cdot, \cdot)$ represents the cosine distance. $F_j^{r,q}$, $F_j^{r,g}$ represents the $j$-th region feature of the query and gallery image, respectively. Analogously, the distance between global features

is given by: $d_g = D(F_g^q, F_g^g)$. Finally, the distance $d$ can be calculated as follows:

$$d = \frac{\sum_{j=1}^N (w_j^q \cdot w_j^g) d_j + d_g}{\sum_{j=1}^N (w_j^q \cdot w_j^g) + 1}. \quad (10)$$

## IV. EXPERIMENTS

### A. Datasets

**Occluded person ReID datasets**. We evaluate the effectiveness of our method in occlusion scenarios. **Occluded-DukeMTMC** [10] contains 15,618 training images, 17,661 gallery images, and 2,210 occluded query images. Occluded-DukeMTMC contains pictures from DukeMTMC-reID [8], while training, query, and gallery contain 9%, 100%, and 10% occluded pictures respectively. **Occluded-REID** [14] is an occluded person ReID dataset obtained from mobile cameras. A total of 2,000 images of 200 individuals are contained within it. For each identity, five full-body images and five occluded images are provided, each with a different viewpoint and an occlusion of a different severity.

**Partial person ReID datasets**. Partial-REID [12] and Partial-iLDIS [11] are two widely used datasets for partial ReID task. **Partial-REID** contains 600 images and 60 people, each of which comprises of five holistic images and five partial images. **Partial-iLIDS** contains 238 images of 119 people from multiple non-overlapping airport cameras, with non-occluded areas manually cropped. In order to make a fair comparison with other competitive methods, we use Market-1501 as a training set and two partial datasets as test sets [10], [22], [27], [28], [48].

**Holistic person ReID datasets**. We also conduct experiments on two widely used holistic person ReID benchmarks Market-1501 [7] and MSMT17 [31] where few images are occluded. In **Market-1501**, the number of training and testing images is 12,936 and 19,732 respectively, with 1,501 identities in all. **MSMT17** dataset contains 126,441 images of 4,101 identities captured from 15 cameras which makes it more challenging. There are 32,621 images of 1,041 identities for training, 93,820 images of 3,060 identities for testing. During inference, 11,659 from 93,820 images are randomly chosen as the query and the other images are viewed as the gallery.

### B. Implementation

**Experimental details:** All training images are resized to $256 \times 128$ and augmented with random horizontal flipping, padding with 10 pixels, random cropping, and random erasing. The batch size is set to 64 with 4 images per person. We use ViT pre-trained on CLIP as our backbone. Adam optimizer is employed with the weight decay factor of 0.0005. The learning rate is initialized as 5e-5 and decreased by a factor of 0.1 at 40th and 70th epochs, respectively, and the training is stopped at 120 epochs. For the prompt training process, we freeze all the other parameters and just train the learnable context vectors under the guidance of $\mathcal{L}_{\mathrm{seg}}$ loss. We employ the Adam optimizer with a learning rate of 5e-5 for this purpose. After 30 epochs of training, we obtain an optimized prompt context, which serves as our learnable prompt specifically designed

TABLE I
PERFORMANCE COMPARISON OF THE OCCLUDED ReID PROBLEM ON THE
OCCLUDED-DUKEMTMC AND OCCLUDED-ReID. THESE SOTA
METHODS ARE CATEGORIZED INTO FOUR GROUPS FROM TOP TO BOTTOM:
HOLISTIC ReID BASED, MATCHING BASED, EXTERNAL CUES BASED, AND
TRANSFORMER BASED.

| Method | Occluded-Duke | | Occluded-REID | |
|---|---|---|---|---|
| | Rank-1 | mAP | Rank-1 | mAP |
| Part-Aligned [49] | 28.8 | 20.2 | - | - |
| PCB [1] | 42.6 | 33.7 | 41.3 | 38.9 |
| Adver Occluded [50] | 44.5 | 32.2 | - | - |
| DSR [27] | 40.8 | 30.4 | 72.8 | 62.8 |
| SFR [28] | 42.3 | 32.0 | - | - |
| FRR [51] | - | - | 78.3 | 68.0 |
| MoS [25] | 61.0 | 49.2 | - | - |
| PVPM [21] | 47.0 | 37.7 | 70.4 | 61.2 |
| PGFA [10] | 51.4 | 37.3 | - | - |
| HOReID [22] | 55.1 | 43.8 | 80.3 | 70.2 |
| GASM [17] | - | - | 74.5 | 65.6 |
| VAN [52] | 62.2 | 46.3 | - | - |
| OAMN [53] | 62.6 | 46.1 | - | - |
| PGFL-KD [23] | 63.0 | 54.1 | 80.7 | 70.3 |
| PAT [54] | 64.5 | 53.6 | 81.6 | 72.1 |
| DRL-Net [55] | 65.8 | 53.9 | - | - |
| FED [56] | 68.1 | 56.4 | 86.3 | 79.3 |
| MSDPA [20] | 70.4 | 61.7 | 81.9 | 77.5 |
| FRT [57] | 70.7 | 61.3 | 80.4 | 71.0 |
| DPM [58] | 71.4 | 61.8 | 85.5 | 79.7 |
| **RGANet** (Ours) | **71.6** | **62.4** | **86.4** | **80.0** |

TABLE II
PERFORMANCE COMPARISON OF THE PARTIAL ReID PROBLEM ON TWO
PARTIAL DATASETS, PARTIAL-REID [12] AND PARTIAL-ILDIS [11].

| Method | Partial-REID | | Partial-iLIDS | |
|---|---|---|---|---|
| | Rank-1 | Rank-3 | Rank-1 | Rank-3 |
| DSR [27] | 50.7 | 70.0 | 58.8 | 67.2 |
| SFR [28] | 56.9 | 78.5 | 63.9 | 74.8 |
| VPM [48] | 67.7 | 81.9 | 65.5 | 74.8 |
| PGFA [10] | 68.0 | 80.0 | 69.1 | 80.9 |
| PVPM [21] | 78.3 | 89.7 | - | - |
| FPR [51] | 81.0 | - | 68.1 | - |
| PGFL-KD [23] | 85.1 | 90.8 | 74.0 | 86.7 |
| HOReID [22] | 85.3 | 91.0 | 72.6 | 86.4 |
| OAMN [53] | 86.0 | - | 77.3 | - |
| FED [56] | 84.6 | 82.3 | - | - |
| MSDPA [20] | 86.3 | 93.3 | 76.5 | 87.4 |
| **RGANet** (Ours) | **87.2** | **93.5** | **77.0** | **87.6** |

MoS [25]), occluded ReID methods depending on external cues (PVPM [21], PGFA [10], HOReID [22], GASM [17], VAN [52], OAMN [53], and PGFL-KD [23]), and methods based on transformer (PAT [54], FED [56], DRL-Net [55], MSDPA [20], FRT [57], and DPM [58]). Our RGANet outperforms all methods by a large margin, attaining 71.6%/62.4% and 86.4%/80.0% in terms of Rank1/mAP on Occluded-DukeMTMC and Occluded-ReID, respectively. For example, PAT [54] exploit transformer to get long range relationships among images to solve occlusions occasions and achieve appealing performance. Our RGANet still surpasses it by +8.8% mAP and +7.9% mAP on Occluded-DukeMTMC and Occluded-ReID, respectively.

**Evaluation on Partial Person ReID Dataset.** In Partial ReID, raw images are manually cropped through a bounding box to solve the matching problem. As a result, severe distortion, misalignment, and occlusions are inevitable, which increases the difficulty of matching. To further study our proposed RGANet, in TABLE II, we also report the comparison of the results on two partial ReID dataset. As can be seen, our RGANet achieves 87.2%/93.5% and 77.0%/87.6% in terms of Rank-1/Rank-3 on Partial-REID and Partial-iLDIS, respectively. Note that, our RGANet is more practical in real-world scene applications because it does not need any external tools during the inference stage.

**Evaluation on Holistic Person ReID Datasets.** To further validate the effectiveness of our proposed RGANet, we perform experiments on Holistic datasets including Market-1501 and MSMT17. As shown in TABLE III, the most recent state-of-the-art methods of ReID can be classified into four groups, namely global feature based models (VCFL [59], MVPM [60], SFT [61], DMML [62], IANet [63], Circle [64], and Mos [25]), part feature based models (DSR [27], PCB [1], VPM [48], and ISP [26]), external cues based methods (MGCAM [16], Pose-transfer [65], PSE [66], SPReID [15], PGFA [10], AANet [67], HOReID [22], GASM [17], and PGFL-KD [23]), and transformer based models (PAT [54], TransReID [6], DRL-Net [55], FED [56], MSDPA [20], FRT [57], PFD [24], and DPM [58]). Our method achieves state-of-the-art performance with 95.5%/89.8% Rank-1/mAP

for the ReID task. The number of the generated regions $N$, the amplification factor $\gamma$, the length of learnable prompts $K$, and the momentum updating rate $m_u$ are set to 4, 20, 8, and 0.3, respectively. Category names fed into CLIP text encoder are set to "`head`", "`upper body`", "`lower body`", "`foot`". The influences of these hyper-parameters will be investigated in the following ablation studies. We freeze the CLIP text encoder during both training and inference. All the experiments are performed with one Nvidia RTX TITAN GPU card using the PyTorch toolbox[2].

**Evaluation Protocols.** As most previous methods in ReID community, Cumulative Matching Characteristic (CMC) curves and the mean Average Precision (mAP) are utilized as metrics to estimate the algorithm. All the experimental results are conducted under the setting of a single query.

### C. Comparison with State-of-the-Art Methods

We compare the proposed RGANet with current state-of-the-art methods on all the above-mentioned ReID datasets of different scenarios.

**Evaluation on Occluded Person ReID Dataset.** To validate the superiority of the RGANet, we evaluate the RGANet on the Occluded Person ReID datasets and show the comparisons in TABLE I. SOTA methods are divided into four mainstreams: holistic ReID methods without special design for occlusions (Part-Aligned [49], PCB [1] and Adver Occluded [50]), occluded ReID methods utilizing matching based approaches (DSR [27], SFR [28], FPR [51], and

[2] http://pytorch.org

TABLE III
COMPARISON WITH SOTA METHODS OVER MARKET-1501. THE
COMPARED METHODS ARE CATEGORIZED INTO THREE MAINSTREAMS:
GLOBAL FEATURE BASED, PART FEATURE BASED, EXTERNAL CUES
BASED, AND TRANSFORMER BASED.

| Methods | Reference | Market-1501 | |
|---|---|---|---|
| | | Rank-1 | mAP |
| VCFL [59] | ICCV'19 | 89.3 | 74.5 |
| MVPM [60] | ICCV'19 | 91.4 | 80.5 |
| SFT [61] | ICCV'19 | 93.4 | 82.7 |
| DMML [62] | ICCV'19 | 93.5 | 81.6 |
| IANet [63] | CVPR'19 | 94.4 | 83.1 |
| Circle [64] | CVPR'20 | 94.2 | 84.9 |
| MoS [25] | AAAI'21 | 94.7 | 86.8 |
| DSR [27] | CVPR'18 | 83.6 | 64.3 |
| PCB [1] | ECCV'18 | 92.3 | 77.4 |
| VPM [48] | CVPR'19 | 93.0 | 80.8 |
| PCB+RPP [1] | ECCV'18 | 93.8 | 81.6 |
| ISP [26] | ECCV'20 | 95.3 | 88.6 |
| MGCAM [16] | CVPR'18 | 83.8 | 74.3 |
| Pose-transfer [65] | CVPR'18 | 87.7 | 68.9 |
| PSE [66] | CVPR'18 | 87.7 | 69.0 |
| SPReID [15] | CVPR'18 | 92.5 | 81.3 |
| PGFA [10] | ICCV'19 | 91.2 | 76.8 |
| AANet [67] | CVPR'19 | 93.9 | 82.5 |
| HOReID [22] | CVPR'20 | 94.2 | 84.9 |
| GASM [17] | CVPR'20 | 94.2 | 84.9 |
| PGFL-KD [23] | CVPR'20 | 94.2 | 84.9 |
| PAT [54] | CVPR'21 | 94.2 | 84.9 |
| TransReID [6] | ICCV'21 | 95.2 | 88.9 |
| DRL-Net [55] | TMM'22 | 94.2 | 84.9 |
| FED [56] | CVPR'22 | 95.0 | 86.3 |
| MSDPA [20] | ACM MM'22 | 95.4 | 89.5 |
| FRT [57] | TIP'22 | 95.5 | 88.1 |
| PFD [24] | AAAI'22 | 95.5 | 89.7 |
| DPM [58] | ACM MM'22 | 95.5 | 89.7 |
| **RGANet** (Ours) | IEEE TIFS | **95.5** | **89.8** |

TABLE IV
PERFORMANCE COMPARISON WITH STATE-OF-THE-ART MODELS ON
MSMT17.

| Methods | Reference | MSMT17 | |
|---|---|---|---|
| | | Rank-1 | mAP |
| MVPM [60] | ICCV'19 | 71.3 | 46.3 |
| SFT [61] | ICCV'19 | 73.6 | 47.6 |
| IANet [63] | CVPR'19 | 75.5 | 46.8 |
| DG-Net [68] | CVPR'19 | 77.2 | 52.3 |
| OSNet [69] | ICCV'19 | 78.7 | 52.9 |
| CBN [70] | ECCV'20 | 72.8 | 42.9 |
| Circle [64] | CVPR'20 | 76.3 | - |
| SAN [71] | AAAI'20 | 79.2 | 55.7 |
| RGA-SC [3] | CVPR'20 | 80.3 | 57.5 |
| TransReID [6] | ICCV'21 | 86.2 | 69.4 |
| DRL-Net [55] | TMM'22 | 78.4 | 55.3 |
| PFD [24] | AAAI'22 | 83.8 | 64.4 |
| **RGANet** (Ours) | IEEE TIFS | **88.1** | **72.3** |

TABLE V
COMPONENTS ANALYSIS OF THE PROPOSED RGANET ON
OCCLUDED-DUKEMTMC. (DAI, AND IAI ARE THE ABBREVIATION FOR
DISCRIMINATION-AWARE INDICATOR AND INVARIANCE-AWARE
INDICATOR, RESPECTIVELY.)

| Type | Index | RAM | | Rank-1 | mAP |
|---|---|---|---|---|---|
| | | DAI | IAI | | |
| Baseline | 0 | ✗ | ✗ | 66.6 | 57.6 |
| Fixed-stripe | 1 | ✗ | ✗ | 68.9 | 58.8 |
| | 2 | ✓ | ✗ | 69.2 | 59.6 |
| | 3 | ✗ | ✓ | 70.0 | 60.2 |
| | 4 | ✓ | ✓ | 70.2 | 61.0 |
| RGM | 5 | ✗ | ✗ | 70.6 | 60.8 |
| | 6 | ✓ | ✗ | 71.0 | 61.4 |
| | 7 | ✗ | ✓ | 71.1 | 61.6 |
| | 8 | ✓ | ✓ | 71.6 | 62.4 |

on Market-1501.

As can be seen from TABLE IV. On MSMT17, our method surpasses the previous holistic SOTA method TransReID [6] by a large margin. *e.g.*, +1.9% and +2.9% in terms of Rank-1 and mAP, respectively. The above experiments demonstrate that our method is a superior universal method for person ReID, which has strong robustness towards different tasks.

### D. Ablation Study

We conduct comprehensive ablation studies on Occluded-DukeMTMC to explore the effectiveness of the proposed modules of our RGANet.

**Analysis of Proposed Modules.** In TABLE V, we evaluate the effectiveness of the Region Generation Module (RGM) and Region Assessment Module (RAM), where DAI, and IAI are the abbreviation for discrimination-aware indicator and invariance-aware indicator, respectively. The baseline utilizes pure ViT pre-trained with CLIP without Overlapping Patches [6] as a feature extractor and is trained under the conventional ID loss and triplet loss. It can achieve 66.6% Rank-1 and 57.6% mAP owing to the strong generalization ability of CLIP which already beats most of the SOTA methods. Next, several observations can be drawn as follows. Firstly, When a naive fixed-stripe (Index 1) which shares

the same architecture with PCB [1] is appended into the baseline (Index 0), the performance is promoted slightly. This is because the fixed stripes are susceptible to misalignment and noise brought from the background or occlusion. By contrast, when adding a single RGM (Index 5) based on the baseline, the performance is improved by +4.0% and +3.2% in terms of Rank-1 and mAP, which demonstrates that our RGM is capable of focusing more on extracting discriminative regions while avoiding being affected by noisy information. Secondly, no matter whether based on the fixed-stripe or the RGM models, further improvement can be achieved by adding DAI or IAI. This implies the strong generalization ability of the DAI and IAI. Thirdly, DAI and IAI are complementary to each other, and the combination of them can further improve the RGM results by +1.0% and +1.6% in terms of Rank-1 and mAP based on RGM. Finally, when combining the RGM and RAM (both DAI and IAI) together, RGANet achieves the best performance with **71.6%** Rank-1 and **62.4%** mAP.

**Prompt design**. We compare five prompt designs methods for RGM only to avoid the impact of RAM. The comparison results are shown in TABLE VI. The naive "a photo of a [CLASS]." obtain 57.1% mAP and simply inserting the word "person" before "photo" would improve the result. While adding "part of a person" into the naive prompt would hurt
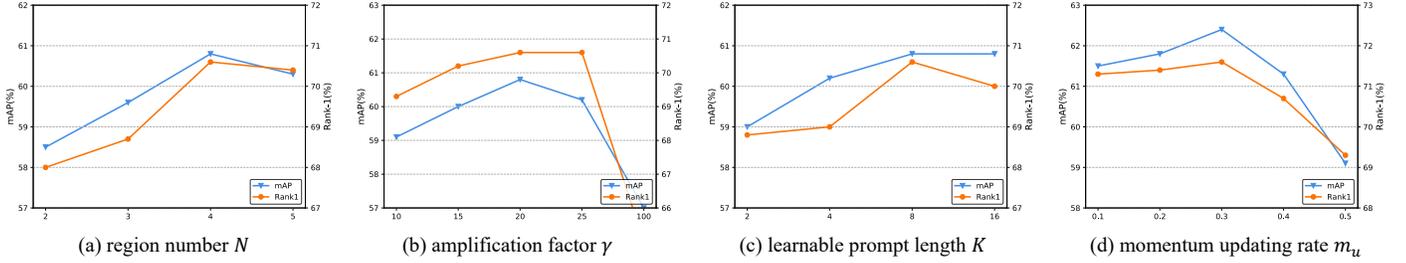
Fig. 3. Effects of three hyper-parameters, the region number $n$, amplification factor $\gamma$, learnable prompt length $K$, and momentum updating rate $m_u$ on Occluded-DukeMTMC.

TABLE VI
PERFORMANCE OF RGANET WITH DIFFERENT PROMPT DESIGNS ON OCCLUDED-DUKE. [CLASS] REPRESENTS THE CLASS TOKEN, AND [LEARNABLE TOKENS] REPRESENT LEARNABLE PROMPTS.

| Prompts | Rank-1 | mAP |
|---|---|---|
| "a photo of a [CLASS]." | 67.6 | 57.1 |
| "a person photo of a [CLASS]." | 68.7 | 58.6 |
| "a photo of a [CLASS] part of a person." | 68.1 | 58.2 |
| "a [CLASS] part of a person." | 69.2 | 59.5 |
| "[Learnable Tokens] + [CLASS]" | 70.6 | 60.8 |

the performance. Therefore, we remove "photo of" and directly utilize "a [CLASS] part of a person.", which improves the result by +2.4% mAP. Finally, the learnable prompt outperforms all the manually searched prompts by +1.3% mAP, clearly showing the generalization power of the learnable prompt. It is worth noting that hand-crafted templates do not need extra training process or external pseudo parsing segmentation and the result surpasses fixed-stripe based methods by +0.7% mAP, which verifies the superiority of our design.

**Evaluation of the Region Number $N$.** To explore the influence of the different number of generated regions, we conduct several experiments based on RGM only without RAM and show the results (Rank-1 and mAP) in Fig. 3 (a). As can be seen, with $N$ increase, the result keeps improving at first and reaches the peak: 70.6% Rank-1 accuracy and 60.8% mAP when $N$ arrives 4. The conclusion is consistent with past methods like PCB and VPM. but overly large region numbers consume more computing resources and slow down the inference speed. Thus, we choose region number 4 to achieve a good balance of accuracy and efficiency. The pre-defined text names for different regions are set to ( "upper body", "lower body", 2 regions in total), ( "head", "upper body", and "lower body", 3 regions in total), ( "head", "upper body", "lower body", and "foot", 4 regions in total), ( "head", "upper body", "lower body", "foot", and "bags", 5 regions in total).

**Analysis of $\gamma$ of RGM.** We conduct experiments on Occluded-DukeMTMC to analyze the impact of $\gamma$ based on RGM only, which denotes the amplification factor for similarity measurement for visual and textual features. As shown in the Fig. 3 (b), too small $\gamma$ leads to poor performance, because the similarity is too small to separate these regions apart. Too large $\gamma$ like 100 may result in unstable training and

loss explosion, which degrades the performance significantly. We set $\gamma$ to 20 in our experiments, which achieves the best result with robustness.

**Impact of $K$ in the learnable prompts.** To investigate the impact of $k$, we perform experiments based on RGM only and show the results in the Fig. 3 (c), our model is not sensitive to the parameter $K$, and the curve is relatively flat. Too small $K$ results in poor performance, because there are not enough context parameters to mine ReID discriminative information. Too large $K$ may result in the redundancy of sample varieties and raise the risk of overfitting, which degrades the performance marginally. Therefore, we set $K$ to 8 in our experiments, which achieves robust performance and keeps efficient.

**Influence of $m_u$.** $m_u$ represents the momentum updating rate used for the invariance-aware indicator in the RAM. As shown in Fig. 3 (d), $m_u = 0.3$ achieves the best performance and is chosen as our default setting. We analyze the performance and conclude that too-small $m_u$ leads to the slow updating of the memory bank feature and too-large $m_u$ results in the drastic updating of the memory bank feature. $m_u = 0.3$ achieves a good balance between the robustness and up-to-dateness of the memory bank.

### E. Visualization of RGANet under Heavy Occlusions

We visualize the generated regions and their assigned confidence values in Fig. 4. For clear explanation, DAI, and IAI are the abbreviations for discrimination-aware indicator, invariance-aware indicator. $\alpha$, $\beta$ and $w$ are discrimination score, invariance score, and final confidence score, respectively. Take Fig. 4 (a) for example, the person's lower body and feet are heavily occluded by the occlusions (car), and the RGM has the capability of detecting the approximate positions of occluded regions, $i.e.$, Region 3 and 4, which are useless for discriminative representation learning. The DAI wrongly assigns Regions 3 and 4 with higher scores, because it can't tell the difference between foreground and background in Regions 3 and 4. The inherent reason is that this man always comes with the car. While the IAI utilizes the auxiliary information from the memory bank to resort to the invariance-level information (complete head, upper body, lower body, and foot). Specifically, Regions 3 and 4 are not invariant features in the memory bank, therefore the IAI values of these regions are lower. $w$ is the mean of the $\alpha$ and the $\beta$, which combines their respective advantages. The experiment results further show
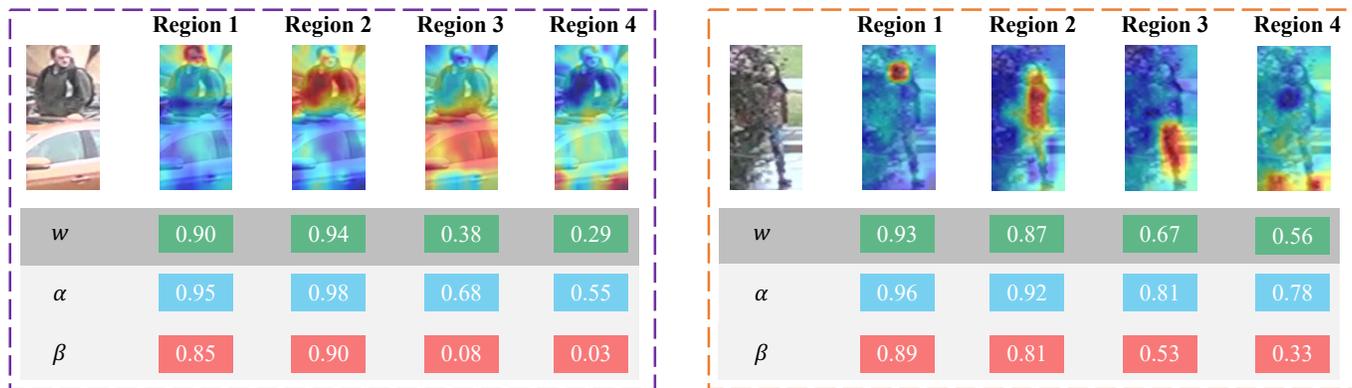
Fig. 4. Visualization of our RGANet under heavy occlusions. $\alpha$, $\beta$ and $w$ are discrimination score, invariance score, and final confidence score, respectively.



Fig. 5. Person image ranking of our baseline and the proposed RGANet are presented on Occluded-DukeMTMC dataset. Green and red rectangles represent right and wrong retrieval results for a given query image, respectively.

that jointly considering the discrimination and the invariance is more effective for person ReID.

### F. Qualitative Retrieval Results

To further show the effectiveness of our proposed RGANet, we provide qualitative retrieval results in Fig. 5. The provided image samples are full of different kinds of occlusions (cars, boards, trees). For each occluded person image from query set, there are the top 10 gallery images that are retrieved by the baseline and our RGANet, respectively. Green and red rectangles represent right and wrong retrieval results for a given query image, respectively. It can be seen that our RGANet is capable of alleviating the occlusions and retrieving the corresponding the same identity images correctly. By contrast, the baseline is susceptible to occlusions and outputs a lot of wrong person images. Therefore, our RGANet can suppress background and occluded parts that are detrimental to the ReID performance.

### G. Computational Complexity

We conduct further experiments in TABLE VII to show that our RGANet not only achieves superior results but also has advantages on inference speed and model complexity. We

TABLE VII
COMPUTATIONAL COMPLEXITY COMPARISON WITH STATE-OF-THE-ART MODELS ON OCCLUDED-DUKE DATASET.

| Methods | FPS | #Params. | Rank-1 | mAP |
|---|---|---|---|---|
| PGFA [10] | 156.4 | 115.4M | 51.4 | 37.3 |
| HOReID [22] | 62.9 | 117.6M | 55.1 | 43.8 |
| PAT [54] | 284.6 | 67.5M | 64.5 | 53.6 |
| DPM [58] | 98.3 | 146.2M | 71.4 | 61.8 |
| **RGANet** (Ours) | 213.9 | 93.5M | 71.6 | 62.4 |

compare recent popular occluded ReID methods: PGFA [10], HOReID [22], and DPM [58]. To ensure a fair comparison, inference batch size is all set to 128. All the experiments are conducted with one Nvidia RTX TITAN GPU using the PyTorch toolbox. The inference time contains both the feature extraction of images and the distance calculation between query and galley images. From the TABLE VII, we can see that our RGANet has relatively small model parameters because we just add several convolution or linear layers apart from our ViT baseline. Besides, the inference speed of RGANet is significantly faster than external tool based methods (PGFA [10] and HOReID [22]) because we get rid of external tools in the inference stage. Although DPM [58]

is also a transformer based method, it uses the operation of Overlapping Patches and multiple $3 \times 3$ convolutional layers to extract mask, which will greatly increase inference time. While PAT has a slight advantage in terms of smaller parameters and faster inference speed, our method significantly outperforms it in terms of performance. The comparisons demonstrate that our proposed RGANet is both compact and efficient.

## V. CONCLUSION

In this paper, to address the challenges of misalignment, background variations, and occlusions in person ReID task, we propose an effective and efficient **R**egion **G**eneration and **A**ssessment **Net**work (**RGANet**), which takes advantage of CLIP to capture the discriminative and invariant region features. Firstly, Region Generation Module (RGM) is designed to automatically search and locate the more discriminative regions. Meanwhile, to obtain reasonable contributions of generated regions, Region Assessment Module (RAM) is proposed to assess each region by a discrimination-aware indicator and an invariance-aware indicator, which are proved to be complementary with each other from both qualitative and quantitative perspectives. Extensive experiments on occluded, partial, and holistic datasets consistently demonstrate the superiority of the proposed approach.

## REFERENCES

[1] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang, "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, pp. 480–496.

[2] G. Wang, Y. Yuan, X. Chen, J. Li, and X. Zhou, "Learning discriminative features with multiple granularities for person re-identification," in *2018 ACM Multimedia Conference on Multimedia Conference*. ACM, 2018, pp. 274–282.

[3] Z. Zhang, C. Lan, W. Zeng, X. Jin, and Z. Chen, "Relation-aware global attention for person re-identification," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[4] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.

[5] X. Jin, C. Lan, W. Zeng, Z. Chen, and L. Zhang, "Style normalization and restitution for generalizable person re-identification," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[6] S. He, H. Luo, P. Wang, F. Wang, H. Li, and W. Jiang, "Transreid: Transformer-based object re-identification," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 15 013–15 022.

[7] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *Computer Vision, IEEE International Conference on Computer Vision*, 2015, pp. 1116–1124.

[8] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[9] Z. Zhong, L. Zheng, D. Cao, and S. Li, "Re-ranking person re-identification with k-reciprocal encoding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1318–1327.

[10] J. Miao, Y. Wu, P. Liu, Y. Ding, and Y. Yang, "Pose-guided feature alignment for occluded person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 542–551.

[11] W.-S. Zheng, S. Gong, and T. Xiang, "Person re-identification by probabilistic relative distance comparison," in *CVPR 2011*. IEEE, 2011, pp. 649–656.

[12] W.-S. Zheng, X. Li, T. Xiang, S. Liao, J. Lai, and S. Gong, "Partial person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4678–4686.

[13] S. Yu, S. Li, D. Chen, R. Zhao, J. Yan, and Y. Qiao, "Cocas: A large-scale clothes changing person dataset for re-identification," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[14] J. Zhuo, Z. Chen, J. Lai, and G. Wang, "Occluded person re-identification," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018, pp. 1–6.

[15] M. M. Kalayeh, E. Basaran, M. Gökmen, M. E. Kamasak, and M. Shah, "Human semantic parsing for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1062–1071.

[16] C. Song, Y. Huang, W. Ouyang, and L. Wang, "Mask-guided contrastive attention model for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1179–1188.

[17] L. He and W. Liu, "Guided saliency feature learning for person re-identification in crowded scenes," in *European Conference on Computer Vision*. Springer, 2020, pp. 357–373.

[18] H. Ding, X. Jiang, B. Shuai, A. Q. Liu, and G. Wang, "Context contrasted feature and gated multi-scale aggregation for scene segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2393–2402.

[19] H. Ding, X. Jiang, A. Q. Liu, N. M. Thalmann, and G. Wang, "Boundary-aware feature propagation for scene segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6819–6829.

[20] X. Cheng, M. Jia, Q. Wang, and J. Zhang, "More is better: Multi-source dynamic parsing attention for occluded person re-identification," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 6840–6849.

[21] S. Gao, J. Wang, H. Lu, and Z. Liu, "Pose-guided visible part matching for occluded person reid," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[22] G. Wang, S. Yang, H. Liu, Z. Wang, Y. Yang, S. Wang, G. Yu, E. Zhou, and J. Sun, "High-order information matters: Learning relation and topology for occluded person re-identification," in *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[23] K. Zheng, C. Lan, W. Zeng, J. Liu, Z. Zhang, and Z.-J. Zha, "Pose-guided feature learning with knowledge distillation for occluded person re-identification," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 4537–4545.

[24] T. Wang, H. Liu, P. Song, T. Guo, and W. Shi, "Pose-guided feature disentangling for occluded person re-identification based on transformer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 2540–2549.

[25] M. Jia, X. Cheng, Y. Zhai, S. Lu, S. Ma, Y. Tian, and J. Zhang, "Matching on sets: Conquer occluded person re-identification without alignment," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1673–1681.

[26] K. Zhu, H. Guo, Z. Liu, M. Tang, and J. Wang, "Identity-guided human semantic parsing for person re-identification," in *European Conference on Computer Vision*. Springer, 2020, pp. 346–363.

[27] L. He, J. Liang, H. Li, and Z. Sun, "Deep spatial feature reconstruction for partial person re-identification: Alignment-free approach," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7073–7082.

[28] L. He, Z. Sun, Y. Zhu, and Y. Wang, "Recognizing partial biometric patterns," *arXiv preprint arXiv:1810.07399*, 2018.

[29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.

[30] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4904–4916.

[31] L. Wei, S. Zhang, W. Gao, and Q. Tian, "Person transfer gan to bridge domain gap for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 79–88.

[32] H. Ding, H. Zhang, C. Liu, and X. Jiang, "Deep interactive image matting with feature propagation," *IEEE Transactions on Image Processing*, vol. 31, pp. 2421–2432, 2022.

[33] J. Mei, Z. Wu, X. Chen, Y. Qiao, H. Ding, and X. Jiang, "Deepdeblur: text image recovery from blur to sharp," *Multimedia tools and applications*, vol. 78, pp. 18 869–18 885, 2019.

[34] D. Li, X. Chen, Z. Zhang, and K. Huang, "Learning deep context-aware features over body and latent parts for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 384–393.

[35] H. Yao, S. Zhang, R. Hong, Y. Zhang, C. Xu, and Q. Tian, "Deep representation learning with part loss for person re-identification," *IEEE Transactions on Image Processing*, vol. 28, no. 6, pp. 2860–2871, 2019.

[36] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-end learning of visual representations from uncurated instructional videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9879–9889.

[37] Z. Wang, J. Yu, A. W. Yu, Z. Dai, Y. Tsvetkov, and Y. Cao, "Simvlm: Simple visual language model pretraining with weak supervision," *arXiv preprint arXiv:2108.10904*, 2021.

[38] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, "Howto100m: Learning a text-video embedding by watching hundred million narrated video clips," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2630–2640.

[39] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.

[40] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Conditional prompt learning for vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 816–16 825.

[41] S. He, H. Ding, and W. Jiang, "Primitive generation and semantic-related alignment for universal zero-shot segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 11 238–11 247.

[42] H. Zhang and H. Ding, "Prototypical matching and open set rejection for zero-shot semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6974–6983.

[43] S. He, H. Ding, and W. Jiang, "Semantic-promoted debiasing and background disambiguation for zero-shot instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 498–19 507.

[44] C. Liu, H. Ding, and X. Jiang, "GRES: Generalized referring expression segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 23 592–23 601.

[45] H. Ding, C. Liu, S. Wang, and X. Jiang, "VLT: vision-language transformer and query generation for referring segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 6, pp. 7900–7916, 2023.

[46] Z. Wang, Z. Fang, J. Wang, and Y. Yang, "Vitaa: Visual-textual attributes alignment in person search by natural language," in *European Conference on Computer Vision*. Springer, 2020, pp. 402–420.

[47] Y. Ge, F. Zhu, D. Chen, R. Zhao *et al.*, "Self-paced contrastive learning with hybrid memory for domain adaptive object re-id," *Advances in Neural Information Processing Systems*, vol. 33, pp. 11 309–11 321, 2020.

[48] Y. Sun, Q. Xu, Y. Li, C. Zhang, Y. Li, S. Wang, and J. Sun, "Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[49] L. Zhao, X. Li, Y. Zhuang, and J. Wang, "Deeply-learned part-aligned representations for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3219–3228.

[50] H. Huang, D. Li, Z. Zhang, X. Chen, and K. Huang, "Adversarially occluded samples for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5098–5107.

[51] L. He, Y. Wang, W. Liu, H. Zhao, Z. Sun, and J. Feng, "Foreground-aware pyramid reconstruction for alignment-free occluded person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8450–8459.

[52] J. Yang, J. Zhang, F. Yu, X. Jiang, M. Zhang, X. Sun, Y.-C. Chen, and W.-S. Zheng, "Learning to know where to see: a visibility-aware approach for occluded person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 11 885–11 894.

[53] P. Chen, W. Liu, P. Dai, J. Liu, Q. Ye, M. Xu, Q. Chen, and R. Ji, "Occlude them all: Occlusion-aware attention network for occluded person re-id," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 11 833–11 842.

[54] Y. Li, J. He, T. Zhang, X. Liu, Y. Zhang, and F. Wu, "Diverse part discovery: Occluded person re-identification with part-aware transformer,"

in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2898–2907.

[55] M. Jia, X. Cheng, S. Lu, and J. Zhang, "Learning disentangled representation implicitly via transformer for occluded person re-identification," *IEEE Transactions on Multimedia*, 2022.

[56] Z. Wang, F. Zhu, S. Tang, R. Zhao, L. He, and J. Song, "Feature erasing and diffusion network for occluded person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4754–4763.

[57] B. Xu, L. He, J. Liang, and Z. Sun, "Learning feature recovery transformer for occluded person re-identification," *IEEE Transactions on Image Processing*, vol. 31, pp. 4651–4662, 2022.

[58] L. Tan, P. Dai, R. Ji, and Y. Wu, "Dynamic prototype mask for occluded person re-identification," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 531–540.

[59] F. Liu and L. Zhang, "View confusion feature learning for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 6639–6648.

[60] H. Sun, Z. Chen, S. Yan, and L. Xu, "Mvp matching: A maximum-value perfect matching for mining hard samples, with application to person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 6737–6747.

[61] C. Luo, Y. Chen, N. Wang, and Z. Zhang, "Spectral feature transformation for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 4976–4985.

[62] G. Chen, T. Zhang, J. Lu, and J. Zhou, "Deep meta metric learning," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 9547–9556.

[63] R. Hou, B. Ma, H. Chang, X. Gu, S. Shan, and X. Chen, "Interaction-and-aggregation network for person re-identification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[64] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, "Circle loss: A unified perspective of pair similarity optimization," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[65] J. Liu, B. Ni, Y. Yan, P. Zhou, S. Cheng, and J. Hu, "Pose transferrable person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4099–4108.

[66] M. Saquib Sarfraz, A. Schumann, A. Eberle, and R. Stiefelhagen, "A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 420–429.

[67] C.-P. Tay, S. Roy, and K.-H. Yap, "Aanet: Attribute attention network for person re-identifications," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[68] Z. Zheng, X. Yang, Z. Yu, L. Zheng, Y. Yang, and J. Kautz, "Joint discriminative and generative learning for person re-identification," in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2138–2147.

[69] K. Zhou, Y. Yang, A. Cavallaro, and T. Xiang, "Omni-scale feature learning for person re-identification," in *ICCV*, 2019, pp. 3702–3712.

[70] Z. Zhuang, L. Wei, L. Xie, T. Zhang, H. Zhang, H. Wu, H. Ai, and Q. Tian, "Rethinking the distribution gap of person re-identification with camera-based batch normalization," in *ECCV*. Springer, 2020, pp. 140–157.

[71] X. Jin, C. Lan, W. Zeng, G. Wei, and Z. Chen, "Semantics-aligned representation learning for person re-identification," in *Thirty-Forth AAAI Conference on Artificial Intelligence*, 2020.