

Dynamic Prototype Mask for Occluded Person Re-Identification

Lei Tan¹, Pingyang Dai^{1*}, Rongrong Ji^{1,2}, Yongjian Wu³

tanlei@stu.xmu.edu.cn, {pydai, rrji}@xmu.edu.cn, littlekenwu@tencent.com

¹Media Analytics and Computing Lab, Department of Artificial Intelligence, School of Informatics, Xiamen University, 361005, China, ²Institute of Artificial Intelligence, Xiamen University, China

³Tencent Youtu Lab, Shanghai, China

ABSTRACT

Although person re-identification has achieved an impressive improvement in recent years, the common occlusion case caused by different obstacles is still an unsettled issue in real application scenarios. Existing methods mainly address this issue by employing body clues provided by an extra network to distinguish the visible part. Nevertheless, the inevitable domain gap between the assistant model and the ReID datasets has highly increased the difficulty to obtain an effective and efficient model. To escape from the extra pre-trained networks and achieve an automatic alignment in an end-to-end trainable network, we propose a novel Dynamic Prototype Mask (DPM) based on two self-evident prior knowledge. Specifically, we first devise a Hierarchical Mask Generator which utilizes the hierarchical semantic to select the visible pattern space between the high-quality holistic prototype and the feature representation of the occluded input image. Under this condition, the occluded representation could be well aligned in a selected subspace spontaneously. Then, to enrich the feature representation of the high-quality holistic prototype and provide a more complete feature space, we introduce a Head Enrich Module to encourage different heads to aggregate different patterns representation in the whole image. Extensive experimental evaluations conducted on occluded and holistic person re-identification benchmarks demonstrate the superior performance of the DPM over the state-of-the-art methods. The code is released at <https://github.com/stone96123/DPM>.

CCS CONCEPTS

• Computing methodologies → Object identification.

KEYWORDS

Occluded Person Re-Identification, Dynamic Prototype Mask

ACM Reference Format:

Lei Tan, Pingyang Dai, Rongrong Ji, Yongjian Wu. 2022. Dynamic Prototype Mask for Occluded Person Re-Identification. In *Proceedings of the 30th ACM*

* Corresponding Author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '22, October 10–14, 2022, Lisboa, Portugal

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9203-7/22/10...\$15.00

<https://doi.org/10.1145/3503161.3547764>

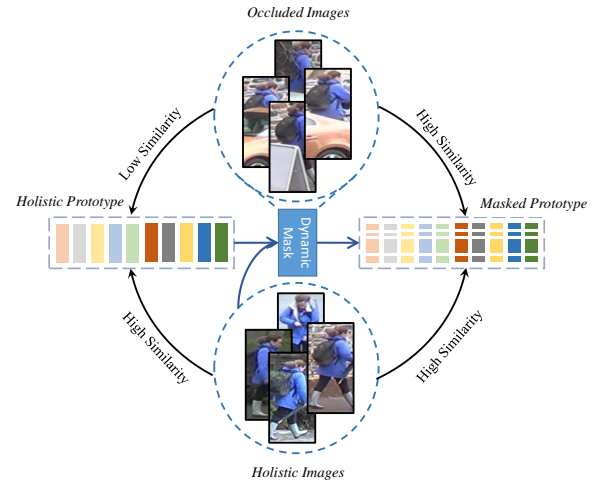


Figure 1: The motivation for Dynamic Prototype Mask (DPM). Based on two prior knowledge (detail discussed in Section 1), we consider the prototype as an ideal holistic representation for each class. Since the occluded image can not provide a complete representation, under DPM, the problem of alignment will be changed to find a visible pattern subspace from the holistic prototype.

International Conference on Multimedia (MM '22), October 10–14, 2022, Lisboa, Portugal. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3503161.3547764>

1 INTRODUCTION

Person re-identification (ReID), aiming to address the problem of matching people over a distributed set of non-overlapping cameras, has attracted intensive attention in the last few years due to its wide applications in surveillance systems [9, 39, 41, 43]. While recent large-scale re-id datasets [36, 47] have provided an ability for deep neural networks to produce a satisfying retrieval performance upon the holistic pedestrian regions, widespread occlusion caused by different obstacles is still an unsettled issue in real application scenarios. This condition in practice inspires a large amount of research effort to explore the occluded person re-identification.

Compared with the general person re-identification problem which assumes the whole body part is available, the main challenge of occluded person re-identification is two folds: Firstly, with the obstacles which cover those discriminative body regions, the valid information in the final feature representation will high-decreased. Even if those valid regions provide valuable information, the final

representation of the whole image may not be corrected at these puny efforts. Secondly, exploring a fine-grid feature representation has been demonstrated as an efficient strategy to achieve an advanced ReID framework [33, 35, 43]. However, occluded person images usually lack several important parts due to the obstacles. Under this condition, those invalid noises will easily provide an ill similarity with similar obstacles and induce an error result. To this end, two typical frameworks have been proposed to tackle the above issues. One of the mainstream frameworks [16, 34] aims to aggregate the information from the whole image and handles the above issue by compensating the invisible body regions by its visible near-neighbor. With the assistance of well-trained human parsing or body key point estimation networks, these methods can easily conduct a topology graph based on the body key point. By passing the information from visible node to invisible node, the influence of occluded regions will largely be alleviated. Although retained the information from the whole image, not only the information from visible near-neighbor may not be so convincing enough, but also the inevitable domain gap between the assistant pre-trained model and the ReID datasets highly increases the difficulty to obtain an effective and efficient model. Alongside the above strategy, discovering and aligning the fine-grid visible body part in the spatial level of the occluded person image are a prevailing and straightforward strategy that has received much attention [10, 27, 30]. By ignoring those invisible parts, these works achieve a significant improvement and visualization result. Nevertheless, to better distinguish the visible/invisible regions, most of the methods in this body of work also rely on extra pre-trained networks to provide body clues and suffer the same domain gap. To make matters worse, those error segmentation or key point results will make the valid nuances be abandoned easily.

Therefore, in this paper, we propose a Dynamic Prototype Mask (DPM) which not only escapes from the extra pre-trained networks but also simultaneously retains the information from the whole image and achieves alignment. The DPM is conducted under two self-evident prior knowledge: 1) During the training, the loss scale of high-quality images which suffer less occlusion will be lower than those of high-occluded images. Since the fully connected layer used for classification could be considered as the bank of prototype for each class, under this perspective, the lower loss scale can be seen as a high similarity between the high-quality sample and its corresponding prototype. In other words, those prototypes for each identity can be considered as a high-quality and complete feature representation that suffers little occlusion. 2) Each channel in the feature representation can be regarded as a response to a specific pattern. This phenomenon in CNN has already been explored by several previous works [2, 17]. For the transformer, the multi-head self-attention directly aggregates the feature based on the similarity from patch to patch. These two prior knowledge indicate that the alignment and matching in the training period for occlusion person re-identification can largely be addressed by selecting a visible pattern subspace for both the input image and its holistic prototype. Motivated by this observation, we introduce the DPM. Different from the spatial attention strategy [3] which takes effect on the feature of the input image itself, the DPM aims to learn a dynamic mask to cut the holistic prototype and select the efficient subspace

for matching. Meanwhile, this processing is totally spontaneous and does not rely on any extra network to provide body clues.

To do this, as shown in Figure 2, the DPM starts with a standard ViT [8] which has demonstrated its superior performance in the computer vision tasks before [5, 15, 21, 24]. Since the key idea of DPM is to generate the prototype mask to select the visible subspace for matching, we first introduce a hierarchical mask generator (HMG) to provide a reliable mask feature. The HMG takes the advantage of the convolutional neural network and aims to evaluate the weight for each channel by the correlation of local information. Meanwhile, we observe that with the network going deeper, the feature representation of each patch will be smoothed and become more similar to each other. Based on high-similar input, it is difficult for pure HMG to provide an efficient prototype mask. Therefore, we add a hierarchical structure to enhance the diversity of the input feature by shallow layers with high diversity. To fully explore the potential of DPM, an Head Enrich Module (HEM) is devised to enrich the feature representation. Specifically, each head in the final transformer block will be encouraged to aggregate different patterns in the whole image. Finally, to evaluate the effectiveness of the proposed DPM, we conduct a series of experiments on both occluded and holistic ReID benchmarks.

The main contributions of the paper are summarized as follows:

- A novel end-to-end trainable network DPM is proposed. DPM not only escapes from the extra pre-trained networks but also simultaneously retains the information from the whole image and achieves automatic alignment.
- To fully explore the potential of DPM, a Hierarchical Mask Generator (HMG) together with a Head Enrich Module (HEM) is introduced. The HMG provides a high-quality sub-space mask via hierarchical semantic information, while the HEM enriches the holistic prototype via diverse heads.
- Extensive experiments on two publicly occluded datasets Occluded-Duke and Occluded-REID demonstrate the superiority of our DPM.

2 RELATED WORKS

2.1 Holistic Person Re-Identification

Holistic person re-identification aims to address the problem of matching people over a distributed set of non-overlapping cameras. The prior works mainly focus on exploring the hand-craft descriptors [22, 26, 38] with a well-designed metric learning strategy [19, 45]. With the resurgence of deep learning, deep feature representation learning has dominated the vision tasks [4, 12, 28, 29, 42]. Luo *et al.* [25] introduce the BN-Neck structure in the CNN-based ReID framework. The research provides a strong baseline for the holistic ReID. Chen *et al.* [1] introduced a high-order attention mechanism to capture and use high-order attention distributions. Zheng *et al.* [46] integrate discriminative and generative learning in a single unified network for person re-identification. Besides using the global feature representation [25, 40, 46], employing part-level features for pedestrian image description to offer fine-grained information is also a mainstream strategy that has been verified as beneficial for person ReID. Methods like PCB [33], MGN [35], and Pyramid [43] horizontally divide the input images or feature maps into several parts to conduct a fine-grid representation. Most recently,

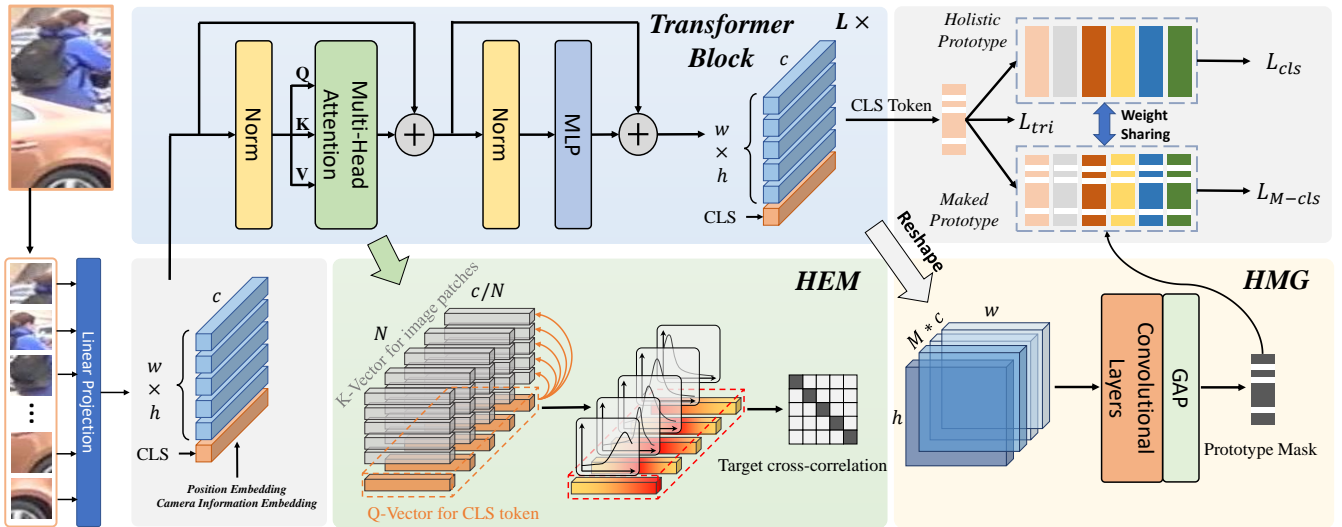


Figure 2: The framework of proposed Dynamic Prototype Mask (DPM). Here, "HEM" denotes the Head Enrich Module, the "HMG" refers to the Hierarchical Mask Generator (HMG). N is the number of the head in transformer block. M is the number of feature maps that be concatenated by HMG. For the input occluded images, HEM encourage the multiple heads in the transformer block to aggregate pattern from different patches. Subsequently, a large amount of enriched feature representation of the dataset will be used to train a holistic prototype for each identity. For each single image with its occluded feature representation, the HMG will provide a dynamic mask to select the appropriate subspace from the holistic prototype.

we have witnessed the thriving of transformer structures from natural language processing to computer vision. TransReID [15] firstly takes the advantage of ViT structure and applies it to the ReID task. Although those methods reach a satisfying performance in the holistic ReID benchmarks, the widely existing occlusion condition is largely ignored. Most of the methods suffer significant performance degradation when being applied to the real-world scenarios which contain the occluded cases.

2.2 Occluded Person Re-Identification

Occluded person re-identification points out the weakness of holistic ReID methods in such occluded cases. The main challenge of occluded ReID lies in the incomplete body information which can not provide high-quality feature representation. To tackle this issue, early works attempt to remove the influence of obstacles in an end-to-end framework and generate the global feature representation from the visible part. Zhuo *et al.* [50] introduce an extra occluded/non-occluded binary classification task to distinguish the occluded images from holistic ones. Chen *et al.* [3] combine an occlusion augmentation scheme with an attention mechanism to precisely capture body parts regardless of the occlusion. Although this kind of work is easy to achieve and shows a good performance before, it always suffer the noise caused by the obstacles which limited the performance upper bound. Therefore, recent methods attempt to avoid such a condition with two typical strategies. The first one aims to aggregate the information from the whole image and handles the above issue by compensating for the invisible body regions by its visible near-neighbor. Wang *et al.* [34] utilize the high-order relation and human-topology information that is based

on keypoint estimation to learn well and robustly aligned features. Hou *et al.* [16] propose a region feature completion module to exploit the long-range spatial contexts from non-occluded regions to predict the features of occluded regions. Though aggregating the information from the visible neighbor node can alleviate the occluded condition, this process is still facing a great challenge when lacking efficient evidence in the neighbor nodes. Meanwhile, using the key point estimate network which is pre-trained on other datasets also faces a challenge to provide reliable results when suffering a domain variation. Another strategy inherits the idea of fine-grid feature representation and aims to match the image between the visible parts. Miao *et al.* [27] introduce the Pose-Guided Feature Alignment (PGFA), exploiting pose landmarks to disentangle visible part information from occlusion noise. Gao *et al.* [10] introduce the Visible Part Matching (PVP) model to learn discriminative part features via a pose-guided attention map. Li *et al.* [21] employ the prototypes to disentangle the fine-grid body part without the help of an extra network in order to achieve satisfying performance. However, most of the fine-grid methods still rely on the extra network to provide body clues and suffer the same domain variation problem. Furthermore, since the fine-grid methods demand strict part prediction to send the feature to its corresponding branch, those incorrect results will make those valuable nuances be ignored easily.

Differing from the above methods, the DPM not only escapes from the extra pre-trained networks but also simultaneously retains the global information from the whole image representation and achieves an automatic alignment.

3 THE PROPOSED METHOD

3.1 Overall Framework

The overview of our proposed DPM framework is illustrated in Figure 2. The DPM adopts a pre-trained ViT [8] to extract the original feature representation from the input images. Herein, we denote the input image as I with the resolution as $H \times W$. We first split the image into $D(h \times w)$ patches with the size after flatten as $x_i \in 1 \times c$. Specifically, it can be described as:

$$D = w \times h = \left\lfloor \frac{H - P + s_d}{s_d} + \frac{W - P + s_d}{s_d} \right\rfloor, \quad (1)$$

where the P and s_d refers to the size of image patch and the step size of sliding window. After the linear projection \mathcal{F} , a learnable class token x_{cls} is attached to aggregate the information from image patches. Before feeding into the transformer block, following the TransReID [15], a learnable position embedding \mathcal{P} and camera embedding \mathcal{C} is added to the patch embeddings to retain positional information and camera information respectively, which can be formulated as:

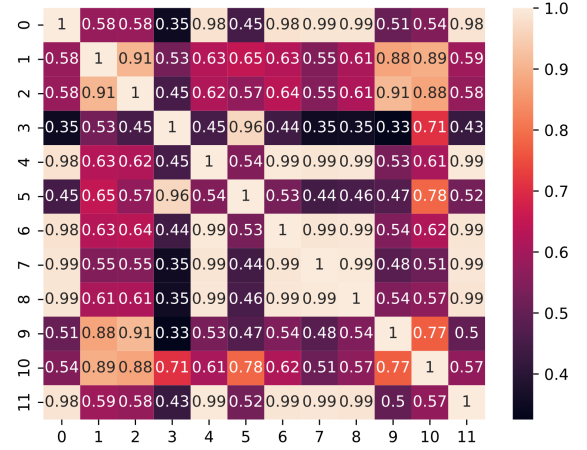
$$z_0 = [x_{cls}^0; \mathcal{F}(x_1^0); \mathcal{F}(x_2^0); \dots; \mathcal{F}(x_D^0)] + \mathcal{P} + \lambda \mathcal{C}, \quad (2)$$

where the z_0 is the input of the transformer blocks. The hyperparameter λ is used to balance the weight of camera embedding. The query and key vectors of the last transformer block are fed into the head enrich module (HEM), which takes the advantage of the multi-attention structure to explore a diverse feature representation for the different heads. Meanwhile, the representation of class-token will be utilized to train a holistic prototype for each class. To well tackle the occluded case, the representation for image patches of the 2_{st} , 4_{th} , 10_{th} , and 12_{th} is concatenated and sent to the hierarchical mask generator (HMG) to provide the dynamic prototype mask for every single input image. Different from spatial attention methods [3] which take effect on the feature map itself, the prototype mask is used to cut the holistic prototype to select the subspace of high-discriminative visible patterns.

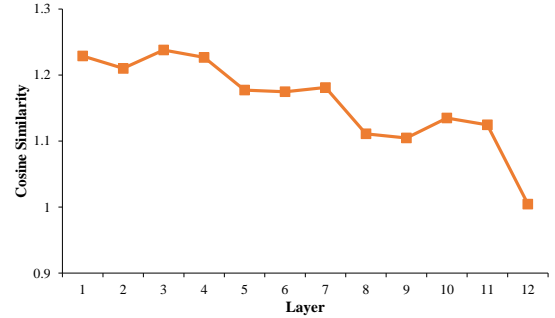
3.2 Hierarchical Mask Generator

Based on the two prior knowledge mentioned before, the main idea of DPM is to explore a spontaneous alignment and select a discriminative subspace for the holistic prototype to match every single image. Therefore, one of the most important parts for achieving the DPM is to generate an efficient prototype mask. In most cases, a pattern within an image is usually conducted by pixels that are spatially concentrated and form a connected component. Therefore, using a local region-based sliding window to weigh the importance of patterns is quite suitable in such an application. In order to provide a high-quality prototype mask, as shown in Figure 2, we apply a convolutional-based mask generator that can take the neighbor nodes into consideration for each patch.

In specific, after the l_{th} block, HMG adopts the reshaped image representation $f_l \in \mathbb{R}^{h \times w \times c}$ by excluding the class-token in the feature representation z_l . Although directly using the image representation provided by the last block seems like the most intuitive strategy, as shown in Figure 3 (b), after calculating the cosine similarity between the most dissimilar image patches in every block, we observe that by passing information through the similarity among



(a) Similarity between multiple heads



(b) Maximum similarity gap between image patches

Figure 3: The motivation for Hierarchical Mask Generator (HMG) and Head Enrich Module (HEM). (a) The cross-correlation matrix between multiple heads' attention maps in the last transformer block. Several class tokens aggregate features from similar image patches which limit the diversity of representation. (b) The maximum similarity gap in each block. Herein, we have ignored the class token. Clearly, the feature of image patches has been smoothed and become more similar to each other with the network going deeper.

image patches, the feature representations will be smoothed and become more similar. Those representations with few discriminative make it hard to provide an efficient prototype mask. To this end, a hierarchical structure is utilized to aggregate those image representations from the shallow layers. Inspired by the success of Swin-Transformer [23] which merges the patches after the 2_{nd} , 4_{th} , 10_{th} transformer block, we also pick the image representations from these three blocks and combine them with the last block as the input for the HMG. The final prototype mask is generated as:

$$M_p = \sigma(\text{Avgpool}(\mathcal{G}(G * (f_1, f_2, \dots, f_L))))), \quad (3)$$

with $f_l = \text{Reshape}[x_1^l; x_2^l; \dots; x_D^l]$.

Herein, the σ refers to the sigmoid function. The \mathcal{G} denotes the convolutional layers of HMG and the $G \in \mathbb{R}^{1 \times L}$ is a binary gate to choose the input for HMG.

Finally, the masked prototype W_{mp} for i_{th} input image is generated by Hadamard product of the row-extended prototype mask M_p^i and the prototype weight matrix W_p as:

$$W_{mp}^i = W_p \odot M_p^i. \quad (4)$$

3.3 Head Enrich Module

Multi-head self-attention which aims to extract difference and discriminative feature representation is considered as the key component to conducting the transformer block. This kind of aggregation of different patterns suits the DPM both in training a complete holistic prototype matrix and selecting an aligned and discriminative subspace. However, as shown in Figure 3 (a), in the original transformer block, there is no explicit optimization to encourage multiple heads to aggregate more nuances in the whole image, which makes it possible for the different heads to have similar feature embedding. The similar head representations will in turn limit the DPM to select a well-aligned subspace.

On account of this condition, we introduce a head enrich module (HEM) to push multiple heads in the class token to obtain diverse patterns in the last transformer block. Generally, the multi-head attention adopts the query matrix Q, key matrix K, and value matrix V to pass the information from different patches. In the training phase, we only employ the class token as the global representation. Therefore, when ignoring the key vector of the class token itself, we can obtain the attention map of $A \in \mathbb{R}^{N \times D}$ between the class token and the image patches as:

$$A(q_L^{cls}, K_L^{img}) = softmax\left(\frac{q_L^{cls} (K_L^{img})^T}{\sqrt{c/N}}\right), \quad (5)$$

where N is the number of heads, q_L^{cls} is the query vector of class token in the L_{th} transformer block, K_L^{img} is the query vector of image patches in the L_{th} transformer block, and N refers to the number of heads. To push the attention map of each head apart, an orthonormal constraint is impose as:

$$\mathcal{L}_{hem} = \left\| \hat{A} \hat{A}^T - \mathbb{I}_N \right\|_F^2, \quad (6)$$

where the $\|\cdot\|_F^2$ is Frobenius norm, the $\mathbb{I}_N \in N \times N$ is the identity matrix, \hat{A} is a normalized A matrix with each row being L2 normalized. With \mathcal{L}_{hem} , the class token could provide richer representation, which not only benefit the learning for holistic prototype and the masked generator.

3.4 Loss Function and Optimization

The DPM contains a two-branch learning framework as the original classification loss and the masked classification loss. Intuitively, we should employ the softmax loss to optimize both two branches. However, this strategy will cause a condition that the mask can not be well learned. After adding the mask, the scale of the loss generated by the masked branch is much smaller than the softmax loss, which can not provide enough power to learn a high-quality prototype mask. Although increasing the weight of the masked

branch may make sense, the softmax loss is not a strong constraint to clustering the samples. Limited constrains also limits the performance of DPM. Inspired by the progress of metric learning [7], we employ an extra angular margin in the original softmax loss to optimize the masked branch. This strategy not only balances the scale of the original branch and masked branch but also highlights the importance of learning a high-quality mask during the training phase. With the extra margin, for input x_{cls}^L with label y^i , the \mathcal{L}_{M-cls} can be given as:

$$\mathcal{L}_{M-cls} = -\frac{1}{B} \sum_{i=1}^B \log \frac{e^{s(\cos(\langle x_{cls}^L, W_{mp}^{y^i} \rangle + m))}}{e^{s(\cos(\langle x_{cls}^L, W_{mp}^{y^i} \rangle + m))} + D_{inter}}, \quad (7)$$

$$\text{with } D_{inter} = \sum_{j=1, j \neq y^i}^C e^{s(\cos(\langle x_{cls}^L, W_{mp}^{y^j} \rangle))},$$

where the B and C refers to the batch size and the number of class. The m denotes the angular margin and s is the hyper-parameter to adjust the scale.

To increase the intra-class similarity and decrease the inter-class similarity, triplet loss \mathcal{L}_{tri} with online hard-mining [31] are combined during the supervise training. Therefore, the overall loss function can be formulated as:

$$\mathcal{L} = \alpha \mathcal{L}_{cls} + (1 - \alpha) \mathcal{L}_{M-cls} + \beta \mathcal{L}_{hem} + \mathcal{L}_{tri} \quad (8)$$

where the α and β is the hyper-parameters to adjust the weight of \mathcal{L}_{M-cls} and \mathcal{L}_{hem} respectively.

4 EXPERIMENT

4.1 Datasets and Experimental Setting

Datasets. To evaluate the effectiveness of the proposed DPM, we conduct extensive experiments on four publicly available ReID benchmarks which include both occluded and holistic person re-identification datasets. The details are as follows.

Occluded-Duke [27] is a large-scale dataset collected from the DukeMTMC for occluded person re-identification. The training set consists of 15,618 images of 702 persons. The testing set contains 2,210 images of 519 persons as the query and 17,661 images of 1,110 persons as the gallery. Until now, Occluded-Duke is still the most challenging dataset for occluded ReID due to its scale.

Occluded-REID [50] is an occluded person dataset captured by mobile cameras. It consists of 2000 images from 200 persons, where each person has 5 whole-body images and 5 occluded person images. Following the evaluation protocol of previous works [3, 10, 34], the Occluded-REID is only used as a testing set. The model used for experiments in this dataset is trained under the training set of Marker-1501 [44].

Market-1501 [44] is a widely-used holistic ReID dataset captured from 6 cameras. It includes 12,936 training images of 751 persons as the training set, 3,368 images of 750 persons as the query, and 19,732 images of 750 persons as the gallery.

DukeMTMC-reID [47] contains 36,441 images of 1,812 persons captured by eight cameras, in which 16,522 images of 702 identities are used as the training set, 2,228 and 16,522 images of 702 persons that do not appear in the training set are used as the query and gallery, respectively.

Method	Occluded-Duke		Occluded-REID	
	R-1	<i>mAP</i>	R-1	<i>mAP</i>
PCB [33]	42.6	33.7	41.3	38.9
Part Bilinear [32]	36.9	-	-	-
FD-GAN [11]	40.8	-	-	-
ISP [49]	62.8	52.3	-	-
TransReID* [15]	66.4	59.2	-	-
DSR [13]	40.8	30.4	72.8	62.8
Ad-Occluded [18]	44.5	32.2	-	-
FPR [14]	-	-	78.3	68.0
PGFA [27]	51.4	37.3	-	-
PVPM+Aug [10]	-	-	70.4	61.2
HOReID [34]	55.1	43.8	80.3	70.2
OAMN [3]	62.6	46.1	-	-
Part-Label [37]	62.2	46.3	81.0	71.0
PAT* [21]	64.5	53.6	81.6	72.1
DPM	71.4	61.8	85.5	79.7

Table 1: Comparison with previous state-of-the-art methods in terms of CMC (%) and *mAP* (%) on Occluded-Duke and Occluded-REID. The symbol * represents methods that employ the transformer structure.

Evaluation Protocol. To verify fair comparison with other methods, we adopt the widely used Cumulative Matching Characteristic (CMC) and mean Average Precision (*mAP*) as evaluation metrics and follow the evaluation settings provided by existing occluded methods [10, 34].

Implementation details. We employ the ViT [8] pre-trained on ImageNet [6] as the backbone network. Particularly, we resize all the input images to 256×128 and adopt commonly used horizontal flipping, padding, random cropping, and random erasing [48] as data augmentation. Following [34, 37], we use extra color jitter augmentation to avoid domain variance when conduct testing in the Occluded-REID. Following the success of TransReID [15], we adopt a lower stride and set λ to 3.0. During the training stage, each mini-batch is conducted by 64 images from 4 identities. In order to strengthen the power of HMG, the training phase is divided into two-step in every iteration. In the first step, we froze the parameter of HMG to train the holistic prototype. In the second step, we froze the parameter except the HMG to train a high-quality prototype mask. During the testing phase, we apply the mask generated by the query image to the gallery images. Then the retrieval stage can still be computed in parallel for each query image. The SGD is utilized as the optimizer, in which the learning rate is initiated as 0.008 with cosine learning rate decay. The hyper-parameters for s and m in Arcface loss are set to 30 and 0.5 respectively in training the Occluded-Duke. In training Occluded-REID, since the strong constraints will induce the overfitting easily when considering the domain variance of Market-1501 and Occluded-REID, we decrease the m to 0.3 and β to 0.01 in training. We implement our DPM with PyTorch and conduct all experiments on a single Nvidia Tesla A100.

Method	Market-1501		DukeMTMC	
	R-1	<i>mA</i>	R-1	<i>mAP</i>
PCB [33]	92.3	71.4	81.8	66.1
MGN [35]	95.7	86.9	88.7	78.4
ISP [49]	95.3	88.6	89.6	80.0
CDNet [20]	95.1	86.0	88.6	76.8
TransReID* [15]	95.2	88.9	90.7	82.0
FPR [14]	95.4	86.6	88.6	78.4
PGFA [27]	91.2	76.8	82.6	65.5
HOReID [34]	94.2	84.9	86.9	75.6
OAMN [3]	93.2	79.8	86.3	72.6
PAT* [21]	95.4	88.0	88.8	78.2
DPM	95.5	89.7	91.0	82.6

Table 2: Comparison with state-of-the-art methods in terms of CMC (%) and *mAP* (%) on Market-1501 and DukeMTMC-reID. The symbol * represents methods that employ the transformer structure.

4.2 Comparison with State-of-the-art Methods

Results on Occluded Datasets. To comprehensively demonstrate the performance of DPM, we evaluate DPM against the previously reported state-of-the-art methods on the Occluded-Duke and Occluded-REID in Table 1, The compared methods include holistic ReID methods [11, 15, 32, 33, 49] and occluded ReID methods [3, 10, 13, 14, 18, 21, 27, 34, 37]. Obviously, the transformer-based structure (PAT, TransReID) has the advantage in solving occluded cases when compared to the convolutional neural network. The DPM also inherits this advantage and further outperforms other transformer-based methods. In the most challenging occluded ReID dataset Occluded-Duke, the DPM reaches an impressive performance, with 71.4% in rank-1 and 61.8% in *mAP*, which at least outperforms other occluded ReID methods with 6.9% and 8.2% in rank-1 and *mAP* respectively. Meanwhile, in the Occluded-REID, the proposed DPM consistently surpasses current state-of-the-art methods. Specifically, the DPM achieves 85.5% in rank-1 accuracy 79.7% in *mAP*, which improves the Rank-1 accuracy by 3.9% and *mAP* by 7.6% over the PAT.

Although not relying on the extra network to provide body clues, the DPM still achieves superior performance in the occluded ReID benchmarks.

Results on Holistic Datasets. Although occluded ReID methods mainly focused on solving the occluded ReID issue, they may suffer a performance decrease in the original holistic ReID task due to incorrect alignment or ignoring of valuable regions. Therefore, in this section, we also evaluate the proposed DPM on the holistic ReID dataset Market-1501 and DukeMTMC-ReID. For better comparison, we select five holistic ReID method [15, 20, 33, 35, 49] and five occluded ReID methods [3, 14, 21, 27, 34].

The results are shown in Figure 2. In the Market-1501 dataset, the DPM gets 95.5% in rank-1 accuracy and 89.7% in *mAP*. In the DukeMTMC-reID, the the DPM gets 91.0% in rank-1 accuracy and

Method	Occluded-Duke			
	R-1	R-5	R-10	mAP
baseline	64.8	80.8	85.8	57.8
+DPM	70.1	82.8	86.9	59.9
+DPM+HR	71.0	82.9	87.2	61.0
+DPM+HR+HEM	71.4	83.7	87.4	61.8

Table 3: Ablation study of each components in DPM on Occluded-Duke. Herein, the HR refers to the hierarchical structure in mask generator, the HEM refers to the head enrich module.

Method	Setting		Occluded-Duke			
	Cls	Mask-Cls	R-1	R-5	R-10	mAP
baseline	\mathcal{S}		64.8	80.8	85.8	57.8
baseline	\mathcal{A}		68.6	81.2	85.9	58.5
DPM	\mathcal{S}	\mathcal{S}	64.6	80.6	85.5	57.3
DPM	\mathcal{A}	\mathcal{A}	68.3	80.0	84.2	57.7
DPM	\mathcal{S}	\mathcal{A}	70.1	82.8	86.9	59.9

Table 4: Performance comparison with different combinations of loss function for classifier and mask-classifier on Occluded-Duke. Here, the \mathcal{S} denotes the branch which is training with the original softmax loss, and the \mathcal{A} denotes the branch which is training with an extra angular margin.

82.6% in mAP. It is clear that DPM shows competitive results in both two holistic ReID datasets when compared to the state-of-the-art holistic ReID methods. When compared to the occluded ReID methods, the DPM outperforms the previous state-of-the-art methods in these two datasets. Overall, the above results show that DPM is a universal framework, which mainly aims to tackle occluded cases. But it will not destroy the performance on the general holistic ReID task.

4.3 Ablation Study

To evaluate the influence of the proposed architectural components. We conduct a series experiments over the occluded-Duke with different settings and show the quantitative results in Table 3. The baseline uses the ViT as backbone and training with the original softmax loss L_{cls} and triplet loss L_{tri} .

Compared to the baseline method, adding the DPM strategy highly improved the performance in both rank-1 accuracy and mAP as 1 and 1 respectively. After adding the hierarchical structure in the mask generator, the performance further increases from 1 and 1 to 1 and 1 in rank-1 and mAP. On the other hand, it also demonstrates that the image representation provided by the last transformer block which lacks sufficient diversity information will limit the performance of mask generator. Meanwhile, HEM also provides significant performance gains as 1 and 1 in rank-1 and mAP based on the above results. The experiments results indicate that all these components have made sense and satisfied their motivation in the



Figure 4: The cross-correlation matrix between multiple heads' attention maps in the class token of the last transformer block after adding the HEM. For better comparison, we visualize the same input image as the Figure 3 (a).

Setting				Occluded-Duke			
P_n	F_n	P	F	R-1	R-5	R-10	mAP
✓				71.4	83.7	87.4	61.8
✓	✓			69.5	82.3	86.8	58.5
		✓		70.2	84.0	87.4	60.8
		✓	✓	69.7	81.6	85.1	58.3

Table 5: Performance comparison with different types of DPM in terms of CMC (%) and mAP (%) on Occluded-Duke. Here, the P and F denotes that the mask is taking effect upon the prototype matrix and the feature representation respectively. The P_n and F_n denote that the mask is taking effect after the L2 normalization.

DPM. All the components contribute to an effective framework consistently and finally result in an impressive performance.

4.4 Discussions

The classification loss function for DPM. As the most important module, how to train a efficient mask generator is one of the main challenges to achieve the the DPM. In Section 3.4, we have mentioned that we utilize an extra angular margin to train the masked branch to avoid inefficient optimization for the masked branch. Therefore, in this part, we conduct a comparison to show the performance of different loss function combinations when training the class branch and masked branch. The results are shown in Table 4. Here, we use the \mathcal{S} to denote the branch which is training with the original softmax loss, and the \mathcal{A} to denote the branch which is training with an extra angular margin. The two baselines are training without the masked branch.

From Table 4, we can observe that adding an extra angular margin could improve the rank-1 accuracy since it can enforce the network to pay more attention to those outlier samples. However, we

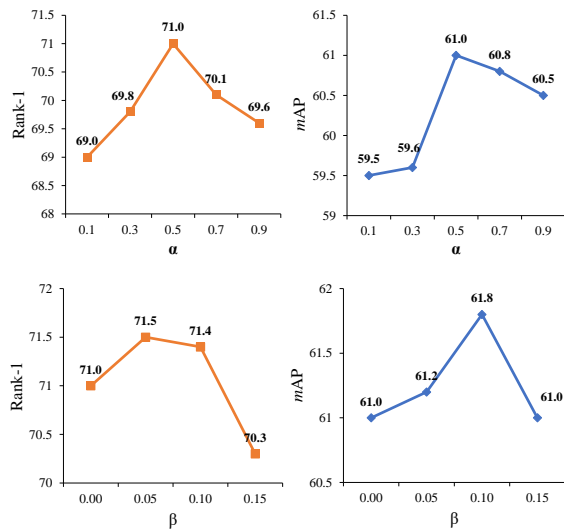


Figure 5: Impact of the hyper-parameters α and β in terms of CMC (%) and mAP (%) on Occluded-Duke. The performance reaches the peak when we set the α to 0.5 and β to 0.1.

can also observe that this strategy can not bring such a significant increment in the mAP , which denotes that the whole distribution may not be ameliorated. Meanwhile, as we have mentioned before, adding an extra mask to select the subspace will further decrease the scale of the loss. It indicates that the mask generator just needs to output the same score for all the channels, the scale loss will still be satisfying. Therefore, as shown in Table 4, when using the same loss function in two branches, the optimization of the mask generator will be limited and the whole network will degrade to optimize the classification branch. After adding an extra angular margin in the masked branch and keeping the original classification branch, the mask generator has been encouraged to adopt a more radical optimization thus alleviating the above dilemma. As shown in the last line in Table 4, this kind of strategy obtain a significant improvement in both rank-1 accuracy and mAP accuracy.

Effect of head enrich module. The HEM aims to enrich the feature representation of multiple heads in the class token to overcome the tendency that more different head has similar representations. The HEM can help not only the training for a holistic prototype with more nuances but also the training for mask generator to select the potential subspace. In Figure 3 (a), we have visualized the cross-correlation matrix between multiple heads' attention map to show the limitation of the original transformer block. Therefore, in this part, we visualized the cross-correlation matrix with the same input once again in Figure 4 after adding the HEM.

By adding the HEM that can explicitly encourage each head to have a different attention map, we could observe that the cross-correlation matrix between multiple heads' attention maps has significantly decreased as shown in Figure 4. It indicates that the class token has received more diverse pattern information from the different image patches. Overall, the visualization well demonstrates the effectiveness of the HEM in such an application.

Impact of the hyper-parameters α and β . As indicated by the loss function in Eq. 8, we set two hyper-parameters α and β to

balance the weight of different components in the overall loss functions. Specifically, the α controls the trade-off between the holistic prototype matrix and the prototype mask, while the β controls the correlation between different heads. Hence, in this part, we conduct empirical experiments to measure the performance of the model under different hyper-parameters settings. When discussing the α , we select the DPM with HMG as the baseline to conduct the experiments. As we have mentioned before the performance is sensitive to the α , a small α will limit the ability to learn a holistic prototype matrix while a large α can not provide enough power to learn a high-quality prototype mask. As shown in the Figure 5, we observe that the performance of rank-1 accuracy and mAP increase linearly when the α is less than 0.5. The performance reaches the peak when the α is set to 0.5 with 71.0% and 61.0% in rank-1 accuracy and mAP . After that, a larger α will decrease the performance.

Based on the model, we further discuss the influence of β . As shown in Figure 5, after adding the HEM into training, the rank-1 accuracy and mAP also get improved when β is less than 0.15. The best performance in rank-1 reaches when the β is set to 0.05 and the best mAP reaches when the β is set to 0.10. Considering the overall performance, we select the β as 0.10 in the further experiments.

Different types of DPM. In DPM, we apply the mask to the holistic prototype matrix, while attention-based strategies explore the spatial attention mask which takes effect on the input image itself to alleviate the noise caused by obstacles. Therefore, in this part, we also make discuss whether the mask should also be applied to the feature representation. Meanwhile, in this part, we also evaluate the influence of L2 normalization in the DPM. Herein, we select the complete DPM as the baseline in the comparison and give an empirical analysis on the Occluded-Duke. The experimental results are shown in Table 5.

Benefits from the great ability of transformer structure that provides an effective representation for the visible part, in both settings, applying the mask upon the feature representation can not provide an extra performance gain. On another side, under both settings, applying the mask on the prototype after the L2 normalization works better than applying the mask before the L2 normalization. Thus, we select the only prototype mask which is applied before the L2 normalization as the final model.

5 CONCLUSION

In this paper, we address the occluded person re-identification with a novel dynamic prototype mask (DPM). The DPM takes the advantage of prototype classification and transfers the alignment in occluded retrieval to the subspace selection task. This strategy not only gets rid of the extra pre-trained networks to provide body clues but also simultaneously retains the information from the global wise and achieves an automatic alignment. Meanwhile, based on the observation in the original DPM framework, we further explore a Hierarchical Mask Generator (HMG) together with a Head Enrich Module (HEM) to fully exploit the potential of DPM. Finally, extensive experiments on occluded and holistic datasets demonstrate the superior performance of DPM.

ACKNOWLEDGMENTS

This work was supported by the National Science Fund for Distinguished Young Scholars (No.62025603), the National Natural Science Foundation of China (No.U1705262, No.62176222, No.62176223, No.62176226, No.62072386, No.62072387, No.62072389, No.62002305, No.61772443, No.61802324 and No.61702136), Guangdong Basic and Applied Basic Research Foundation (No.2019B1515120049), the Natural Science Foundation of Fujian Province of China (No.2021J01002), and the Fundamental Research Funds for the Central Universities (No.20720200077, No.20720200090 and No.20720200091).

REFERENCES

- [1] Binghui Chen, Weihong Deng, and Jiani Hu. 2019. Mixed High-Order Attention Network for Person Re-Identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [2] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. 2017. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5659–5667.
- [3] Peixian Chen, Wenfeng Liu, Pingyang Dai, Jianzhuang Liu, Qixiang Ye, Mingliang Xu, Qi'an Chen, and Rongrong Ji. 2021. Occlude them all: Occlusion-aware attention network for occluded person re-id. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 11833–11842.
- [4] Zhiwei Chen, Liujuan Cao, Yunhang Shen, Feihong Lian, Yongjian Wu, and Rongrong Ji. 2021. E2Net: Excitatory-expansive learning for weakly supervised object Localization. In *Proceedings of the 29th ACM International Conference on Multimedia*. 573–581.
- [5] Zhiwei Chen, Changan Wang, Yabiao Wang, Guannan Jiang, Yunhang Shen, Ying Tai, Chengjie Wang, Wei Zhang, and Liujuan Cao. 2022. Lctr: On awakening the local continuity of transformer for weakly supervised object localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 410–418.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Proceedings of the CVPR*. 248–255.
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4690–4699.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [9] Chanh Eom and Bumsub Ham. 2019. Learning disentangled representation for robust person re-identification. In *Proceedings of the NeurIPS*. 5297–5308.
- [10] Shang Gao, Jingya Wang, Huchuan Lu, and Zimo Liu. 2020. Pose-guided visible part matching for occluded person ReID. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11744–11752.
- [11] Yixiao Ge, Zhuowan Li, Haiyu Zhao, Guojun Yin, Shuai Yi, Xiaogang Wang, et al. 2018. Fd-gan: Pose-guided feature distilling gan for robust person re-identification. *Advances in neural information processing systems* 31 (2018).
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [13] Lingxiao He, Jian Liang, Haiqing Li, and Zhenan Sun. 2018. Deep spatial feature reconstruction for partial person re-identification: Alignment-free approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7073–7082.
- [14] Lingxiao He, Yinggang Wang, Wu Liu, He Zhao, Zhenan Sun, and Jiashi Feng. 2019. Foreground-aware pyramid reconstruction for alignment-free occluded person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 8450–8459.
- [15] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. 2021. Transreid: Transformer-based object re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 15013–15022.
- [16] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. 2021. Feature completion for occluded person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [17] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.
- [18] Houjing Huang, Dangwei Li, Zhang Zhang, Xiaotang Chen, and Kaiqi Huang. 2018. Adversarially occluded samples for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5098–5107.
- [19] Martin Koestinger, Martin Hirzer, Paul Wohlhart, Peter M Roth, and Horst Bischof. 2012. Large scale metric learning from equivalence constraints. In *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2288–2295.
- [20] Hanjun Li, Gaojie Wu, and Wei-Shi Zheng. 2021. Combined depth space based architecture search for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6729–6738.
- [21] Yulin Li, Jianfeng He, Tianzhu Zhang, Xiang Liu, Yongdong Zhang, and Feng Wu. 2021. Diverse part discovery: Occluded person re-identification with part-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2898–2907.
- [22] Shengcai Liao, Yang Hu, Xiangyu Zhu, and Stan Z Li. 2015. Person re-identification by local maximal occurrence representation and metric learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2197–2206.
- [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10012–10022.
- [24] Gen Luo, Yiyi Zhou, Xiaoshuai Sun, Yan Wang, Liujuan Cao, Yongjian Wu, Feiyue Huang, and Rongrong Ji. 2022. Towards Lightweight Transformer Via Group-Wise Transformation for Vision-and-Language Tasks. *IEEE Transactions on Image Processing* 31 (2022), 3386–3398.
- [25] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. 2019. Bag of tricks and a strong baseline for deep person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 0–0.
- [26] Bingpeng Ma, Yu Su, and Frederic Jurie. 2014. Covariance descriptor based on bio-inspired features for person re-identification and face verification. *Image and Vision Computing* 32, 6–7 (2014), 379–390.
- [27] Jiayu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. 2019. Pose-guided feature alignment for occluded person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 542–551.
- [28] Yongqiang Mou, Lei Tan, Hui Yang, Jingying Chen, Leyuan Liu, Rui Yan, and Yaohong Huang. 2020. Plugnet: Degradation aware scene text recognition supervised by a pluggable super-resolution unit. In *European Conference on Computer Vision*. Springer, 158–174.
- [29] Jun Peng, Yiyi Zhou, Xiaoshuai Sun, Liujuan Cao, Yongjian Wu, Feiyue Huang, and Rongrong Ji. 2021. Knowledge-Driven Generative Adversarial Network for Text-to-Image Synthesis. *IEEE Transactions on Multimedia* (2021).
- [30] Xuelin Qian, Yanwei Fu, Tao Xiang, Wenxuan Wang, Jie Qiu, Yang Wu, Yu-Gang Jiang, and Xiangyang Xue. 2018. Pose-normalized image generation for person re-identification. In *Proceedings of the ECCV*. 650–667.
- [31] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the CVPR*. 815–823.
- [32] Yumin Suh, Jingdong Wang, Siyu Tang, Tao Mei, and Kyoung Mu Lee. 2018. Part-aligned bilinear representations for person re-identification. In *Proceedings of the European conference on computer vision (ECCV)*. 402–419.
- [33] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. 2018. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *Proceedings of the ECCV*. 480–496.
- [34] Guan'an Wang, Shuo Yang, Huan Yu Liu, Zhicheng Wang, Yang Yang, Shuliang Wang, Gang Yu, Erjin Zhou, and Jian Sun. 2020. High-order information matters: Learning relation and topology for occluded person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6449–6458.
- [35] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. 2018. Learning discriminative features with multiple granularities for person re-identification. In *Proceedings of the ACM MM*. 274–282.
- [36] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. 2018. Person transfer gan to bridge domain gap for person re-identification. In *Proceedings of the CVPR*. 79–88.
- [37] Jinrui Yang, Jiawei Zhang, Fufu Yu, Xinyang Jiang, Mengdan Zhang, Xing Sun, Ying-Cong Chen, and Wei-Shi Zheng. 2021. 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*. 11885–11894.
- [38] Yang Yang, Jimei Yang, Junjie Yan, Shengcai Liao, Dong Yi, and Stan Z Li. 2014. Salient color names for person re-identification. In *European conference on computer vision*. Springer, 536–551.
- [39] Mang Ye, Weijian Ruan, Bo Du, and Mike Zheng Shou. 2021. Channel Augmented Joint Learning for Visible-Infrared Recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13567–13576.
- [40] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. 2021. Deep learning for person re-identification: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [41] Yunpeng Zhai, Shijian Lu, Qixiang Ye, Xuebo Shan, Jie Chen, Rongrong Ji, and Yonghong Tian. 2020. AD-Cluster: Augmented Discriminative Clustering for Domain Adaptive Person Re-Identification. In *Proceedings of the CVPR*.
- [42] Yukang Zhang, Yan Yan, Yang Lu, and Hanzhi Wang. 2021. Towards a Unified Middle Modality Learning for Visible-Infrared Person Re-Identification. In *Proceedings of the 29th ACM International Conference on Multimedia*. 788–796.

- [43] Feng Zheng, Cheng Deng, Xing Sun, Xinyang Jiang, Xiaowei Guo, Zongqiao Yu, Feiyue Huang, and Rongrong Ji. 2019. Pyramidal Person Re-Identification via Multi-Loss Dynamic Training. In *Proceedings of the CVPR*. 8514–8522.
- [44] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*. 1116–1124.
- [45] Wei-Shi Zheng, Shaogang Gong, and Tao Xiang. 2012. Reidentification by relative distance comparison. *IEEE transactions on pattern analysis and machine intelligence* 35, 3 (2012), 653–668.
- [46] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. 2019. Joint discriminative and generative learning for person re-identification. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2138–2147.
- [47] Zhedong Zheng, Liang Zheng, and Yi Yang. 2017. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE international conference on computer vision*. 3754–3762.
- [48] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. 2020. Random Erasing Data Augmentation. In *Proceedings of the AAAI*.
- [49] Kuan Zhu, Haiyun Guo, Zhiwei Liu, Ming Tang, and Jinqiao Wang. 2020. Identity-guided human semantic parsing for person re-identification. In *European Conference on Computer Vision*. Springer, 346–363.
- [50] Jiaxuan Zhuo, Zeyu Chen, Jianhuang Lai, and Guangcong Wang. 2018. Occluded person re-identification. In *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 1–6.