

Video Person Re-identification using Attribute-enhanced Features

Tianrui Chai, Zhiyuan Chen, Annan Li*, *Member, IEEE*, Jiaxin Chen, Xinyu Mei, and Yunhong Wang, *Fellow, IEEE*

Abstract—Video-based person re-identification (Re-ID) which aims to associate people across non-overlapping cameras using surveillance video is a challenging task. Pedestrian attribute, such as gender, age and clothing characteristics contains rich and supplementary information but is less explored in video person Re-ID. In this work, we propose a novel network architecture named Attribute Saliency Assisted Network (ASA-Net) for attribute-assisted video person Re-ID, which achieved considerable improvement to existing works by two methods. First, to learn a better separation of the target from background, we propose to learn the visual attention from middle-level attribute instead of high-level identities. The proposed Attribute Salient Region Enhance (ASRE) module can attend more accurately on the body of pedestrian. Second, we found that many identity-irrelevant but object or subject-relevant factors like the view angle and movement of the target pedestrian can greatly influence the two dimensional appearance of a pedestrian. This problem can be mitigated by investigating both identity-relevant and identity-irrelevant attributes via a novel triplet loss which is referred as the Pose & Motion-Invariant (PMI) triplet loss. Extensive experiments on MARS and DukeMTMC-VideoReID datasets show that our method outperforms the state-of-the-art methods. Also, the visualizations of learning results well-explain the mechanism of how the improvement is achieved.

Index Terms—video-based person Re-ID, pedestrian attribute, attribute salient region enhance, pose & motion-invariant triplet loss

I. INTRODUCTION

PERSON re-identification (Re-ID) is a specific person retrieval problem which searches a particular person in a query image or video across non-overlapping cameras. It has a wide range of applications, such as person tracking, criminal searching and activity analysis and has attracted more and more attention.

In the past few years, many methods have been proposed for image-based person Re-ID [5], [4], [2], [68], [69], [70] and video-based person Re-ID [39], [38], [36], [72]. Remarkable results have been achieved with them. These methods can be divided into two categories: i.e. feature or representation learning [5], [4], [2], [39], [38], [36] and metric learning [68], [69], [70] respectively. The former focuses on learning discriminative appearance feature of pedestrian, while the latter focuses on deriving appropriate distance metric for matching.

T. Chai, A. Li, J. Chen, X. Mei and Y. Wang are with the School of Computer Science and Engineering, Beihang University, Beijing, China, 100191, e-mail: {trchai, liannan, jiaxinchen, xyme, yhwang}@buaa.edu.cn.

Z. Chen is with the Alibaba group, e-mail: dechen@buaa.edu.cn. This work was done when Chen was a master student at Beihang University.

A. Li is the corresponding author.



Fig. 1: Exemplar sequences and corresponding attribute annotations. As shown in the left image of (a) and (c), it's hard to identify all the identity-relevant attribute from a single image. However, these attributes can be clearly observed from other frames. Despite the identity-relevant attributes, we find identity-irrelevant attributes are also helpful. Similarly from the left single image of (b), it's hard to say whether the subject is walking or running. But with the help of adjacent frames we can infer that he is running. Also, the waiting behavior of the woman in sequence (d) is ambiguous when only one image is available.

No matter the aforementioned feature learning or the metric learning, they are all supervised by the identity which is at a high level of visual semantics. Since prior arts found that middle-level attributes such as gender, age and clothing characteristics can provide additional information, they have been incorporated in person Re-ID. Existing approaches for attribute-enhanced or assisted person Re-ID are mainly based on static image. Some works [28], [29], [76] considered attributes as discriminative features while the others [30], [31], [32], [53], [77] take attributes as a supervision to help co-training.

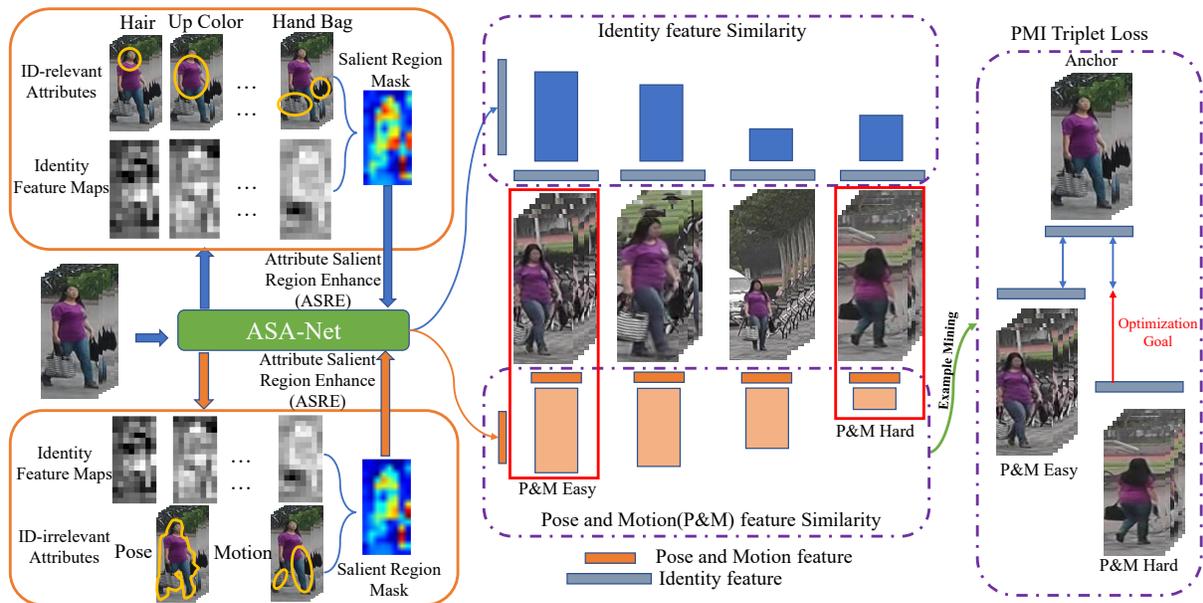


Fig. 2: Overview of our methods. When training the ASA-Net, both ID-relevant and ID-irrelevant attributes are used to enhance features. Pose & motion feature as well as the identity feature are calculated with ASA-Net. And Pose & motion feature is used for mining robust identity feature based on the Pose & Motion-Invariant (PMI) triplet loss.

However, though pedestrian attributes show good assistance in image-based person Re-ID, we argue that attributes haven't fully play a role due to the limitations of static image. As shown in (a) and (c) of Figure 1, attributes such as *shoulder bag*, *shoes* and *down color* cannot be observed from a single frame due to partial occlusions. But these attributes can be clearly identified from other frames of the sequence. A video or an image sequence can provide more complete pedestrian attribute information. Consequently, exploring attribute cues in video is a better way of tackling attribute-assisted person Re-ID. Some efforts have been made to improving video person Re-ID by attributes. Song et al. [33] and Zhao et al. [34] first introduced attributes into video-based person Re-ID. Li et al. [35] designed an encoder-decoder structure to capture the motion information to enhance the Re-ID model. However, such pioneer works have two obvious limitations.

First, although video or image sequence can provide redundant information to person Re-ID, there are still some problems to be solved. Specifically, besides outer factors like background clutter or temporal occlusions, the viewpoint/pose difference, movement variation or other object-relevant but identity-irrelevant (ID-irrelevant) factors can also result in a dramatic change in two dimensional appearance. Existing approaches for video-based person Re-ID seldom include any mechanism of explicitly identifying such factor. Consequently, some ID-irrelevant pattern will be inevitably treated as part of the subject-specific feature, and results in errors of identification.

Second, we observe that one of the fundamental challenge of pedestrian recognition is the template a fact that the template is usually a simple rectangle but the shape or silhouette of a person can be various. Therefore, separating background element from the target pedestrian is always an important issue. In video-based settings, the bounding box of target pedestrian

is obtained by automatic detector and tracker, by which some errors are inevitable. The inaccurate box makes the target-background separation issue worse. To mitigate this problem, explicit segmentation [44], [88] and visual attention [55], [59], [64], [82], [86] have been adopted. The former requires a bulk of pixel-level annotation, which is difficult to obtain, and the latter is usually supervised by the identity labels. We argue that although identity labels can tell who is the people, the clue they provide about the spatial location of the people is coarse-grained. In contrast, attributes can provide various fine-grained information about the figure, motion/action and the status of carrying article of the target people. Therefore, learning from attributes is a better way of tackling the target-background separation issue.

Based on the above observations, we propose a comprehensive architecture or attribute-assisted video person re-identification. It is referred as Attribute Saliency Assisted Network (ASA-Net). Its overview is shown in Figure 2. ASA-Net combines both attribute and identity recognition and it is a five-branch structure: two attribute branches for learning attribute feature, one base branch for extracting original identity features and two branches for learning the attribute-enhanced identity feature. What makes ASA-Net different from existing method lies in two aspects: First, to address the target-background separation issue in an attribute-driven manner, we design the Attribute Salient Region Enhance (ASRE) module. By introducing this module the clutteredness caused by the background elements can be considerably reduced. Second, to deal with the ID-irrelevant but subject-relevant factors, we propose the Pose & Motion-Invariant triplet loss (PMI triplet Loss). The PMI triplet Loss aims to mine the hardest samples through the distance of pose and motion states, and to reduce the intra-class distance caused by pose and motion state differences. Using the attribute annotation provided by

Chen et al. [47] the proposed ASA-Net achieved promising performance on MARS [41] and DukeMTMC-VideoReID [42] datasets, which well demonstrate the effectiveness of our method.

In summary, our contributions are three-fold:

- 1) We design a five-branch network named ASA-Net which integrates video-based attribute recognition and person Re-ID. The novel Attribute Salient Region Enhance (ASRE) module can effectively focus on the attribute salient region and enhance the feature.
- 2) We propose pose & motion-invariant triplet loss (PMI Triplet Loss) which mines the hardest samples with the middle-level ID-irrelevant attributes to reduce the difference introduced by the change of pose and motion.
- 3) We evaluate proposed ASA-Net and PMI triplet loss on MARS [41] and DukeMTMC-VideoReID [42], on which the proposed method outperforms the state-of-the-art. The effectiveness of our method is validated by both experiments and the visualization of learning results.

The rest of the paper is organized as follows: we overview works closely related to our work in Section II. We elaborate the details of the proposed method in Section III, and provide experimental results and analysis of our method by comparing with the state-of-the-art approaches in Section IV. Finally we conclude in Section V.

II. RELATED WORK

In this section, we provide a brief review of relevant existing Re-ID methods. They can be classified into three categories: image-based person Re-ID, video-based person Re-ID, and attribute-assisted person Re-ID in both image-based and video-based settings.

A. Image-based Person Re-ID

Person Re-ID is a challenging task which has been studied for years. Existing works on image-based person Re-ID can be divided into two categories: the first one focus on learning a discriminative ID-level representation for the input pedestrian image [5], [4], [2], [1], [3], [86] while the second one focus on learning robust distance metrics [6], [7], [81], [84]. In the early years, methods focused on the design and extraction of hand-crafted features [9], [10]. With the great success of deep learning [11], convolutional neural networks have been adopted in many approaches [13], [14], [12], [15] and have achieved great improvements. Although achieved good results, image-based settings actually avoid the negative influence of inaccurate pedestrian detection and tracking. To meet the requirements from real-world application, video-based person Re-ID which is based on automatic detection and tracking is becoming the mainstream.

B. Video-based Person Re-ID

Compared with image, temporal relation between frames in a pedestrian sequence or tracklet can make up for important identity information lost by single image due to occlusion,

view change and other factors. The same as image-based person Re-ID, traditional approaches for video-based person re-id also focus on finding efficient hand-crafted descriptors [16], [17]. In the past years, deep learning based approaches [18], [19], [20], [21] for video-based person Re-ID have been proposed and show good performance. Recently, some work focus on the mining of spatial-temporal relationships. Yang et al. [36] and Yan et al. [40] construct nodes according to the body part and the spatial-temporal features are aggregated by using graph convolution network. Gu et al. [39] put his hand to the problem of spatial misalignment of video clips and proposed the AP3D net to align the adjacent in pixel level. Zhang et al. [43] proposed the MGRA net to reduce redundancy within a clip.

C. Attribute-assisted Person Re-ID

Pedestrian attribute recognition [22], [23], [24], [25] has been a hot topic because of its practicability and great demand from video surveillance. Due to the natural similarity between pedestrian attribute learning and person Re-ID, it is natural to take them into consideration at the same time. Compared with the sole person Re-ID, attribute can provide a higher-level semantic discriminant information such as gender, ethnic and age. Some early methods use the attributes to increase the discrimination of identity feature [26], [27], [28], [29], [78], while some other methods consider attribute as a supervision to help co-training [30], [31], [32], [53], [79]. Also, attributes can be considered as a metric for identification [80].

Song et al. [33] first applied attribute information into video-based person Re-ID. Zhao et al. [34] transferred the knowledge of attribute dataset into video-based person Re-ID. Li [35] first introduced the concept of motion into the attributes and used attributes to enhance the video-based person Re-ID. Although these works show the effectiveness of attribute in video person Re-ID, the lack of attribute annotation limits the performance improvement and model development. Also, the attributes considered by these work are consistent with those in image-based person Re-ID, and the unique pedestrian ID-irrelevant attributes which have a great influence on the appearance are overlook. To this end, we propose a comprehensive study on attribute-assisted video person re-identification.

III. METHOD

As shown in Fig. 3, the framework for our proposed model consists of three branches. The upper branch is the ID-relevant attribute branch which extracting features of ID-relevant attributes (gender, shoes, hair, up color...) and make predictions of these attributes. The middle branch is the original branch which extracting appearance feature of the pedestrian for identity recognition. The bottom branch is the ID-irrelevant attribute branch. It is similar to the upper branch but the attributes are replaced with those unrelated to identity such as motion and pose. Every branch is not an isolated island. We use the proposed ASRE module to enhance the feature map extracted from the original branch on the spatial level according to the feature map extracted from the attribute

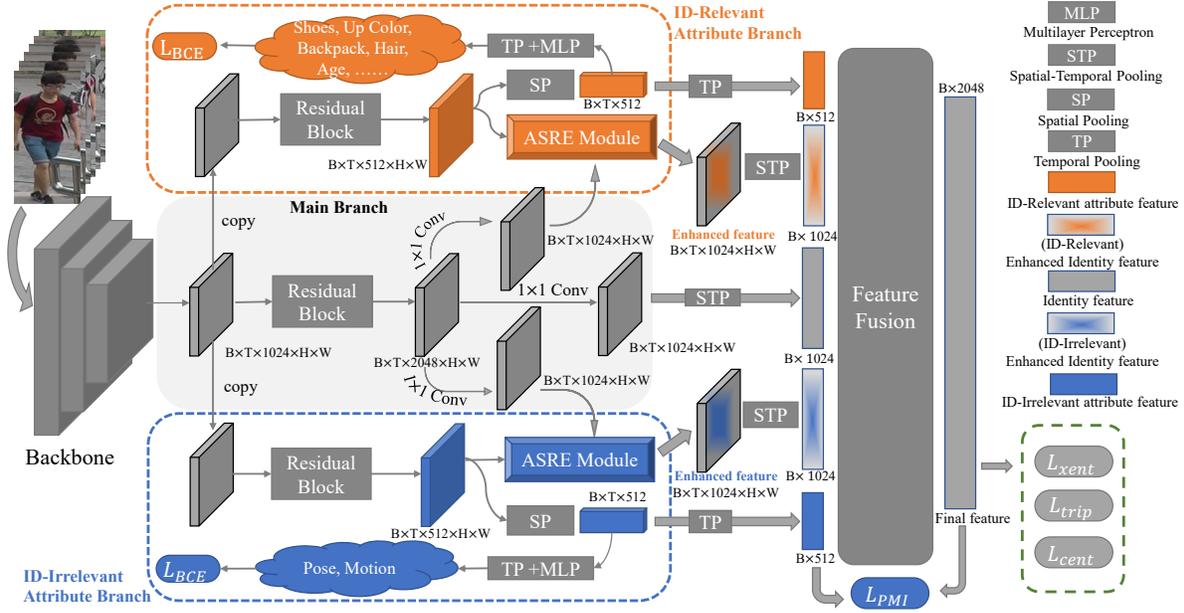


Fig. 3: Overall framework of our proposed ASA-Net. The input sequence has T frames and we use a CNN backbone to extract the feature map of each frame. Then we send the feature map into three branches. The middle row is the identity branch which extracts original identity features. The ID-relevant attribute branch extracts the features for recognizing attributes like gender, clothes, hair and etc. The identity features is enhanced by attribute feature using the Attribute-Salient-Region-enhance (ASRE) module. The framework of ID-irrelevant attribute branch is similar to the ID-relevant branch. Finally we fuse the obtained five features to get the final feature. At the same time, we also construct the Pose&Motion-invariant (PMI) loss to eliminate the influence caused by pose and motion according to the ID-irrelevant attribute features.

branch. Five branches of features will be obtained with ASA-Net. Especially, the feature extracted from ID-irrelevant branch are used to calculate the proposed PMI triplet loss.

In this section, we will first introduce the overall framework of our ASA-Net. The next we will introduce the ASRE module, feature fusion strategy, the PMI triplet loss and model optimization in detail respectively.

A. Overall Framework

Given a pedestrian video clip or tracklet, we denote it as $I = \{I_1, I_2, I_3, \dots, I_T\}$ where T is the number of frames sampled from the clip. The feature map extracted by the backbone network is $X_{origin} = [X_1, X_2, \dots, X_T]$ where $X_i \in \mathbb{R}^{C \times H \times W}$ and $X_{origin} \in \mathbb{R}^{T \times C \times H \times W}$. We copy the feature map X_{origin} twice and send them into attribute branches. After their own residual module [11] of the three branches, we can get three basic feature maps X_{re_attr} , X_{ID} , X_{ir_attr} . The process can be expressed as:

$$\begin{aligned}
 X_{origin} &= F_{backbone}(I), \\
 X_{re_attr} &= F_{residual_re}(X_{origin}) \\
 X_{ID} &= F_{residual_ID}(X_{origin}) \\
 X_{ir_attr} &= F_{residual_ir}(X_{origin}),
 \end{aligned} \tag{1}$$

where $X_{re_attr}, X_{ir_attr} \in \mathbb{R}^{T \times \frac{C}{2} \times H \times W}$ and $X_{ID} \in \mathbb{R}^{T \times 2C \times H \times W}$. $F_{residual_re}$, $F_{residual_ID}$ and $F_{residual_ir}$ are residual module [11] of the corresponding dimension.

In the ID-relevant attribute branch, X_{re_attr} are used to calculate the final ID-relevant attribute feature f_{re_attr} which

is used for ID-relevant attribute prediction and the final identity feature calculation. At the same time, the enhanced feature map X_{re_ID} is calculated with the ASRE module according to the ID-relevant attribute feature map X_{re_attr} and the identity feature map X_{ID} . Finally, the ID-relevant attribute enhanced feature f_{re_ID} are got according to the X_{re_ID} with the spatial-temporal pooling. The ID-relevant attribute branch is formulated as following:

$$\begin{aligned}
 X_{re_ID} &= ASRE(X_{re_attr}, F_{re}(X_{ID})), \\
 f_{re_ID} &= W_{re}(F_{TP}(F_{SP}(X_{re_ID}))),
 \end{aligned} \tag{2}$$

where $ASRE$, F_{TP} , F_{SP} , W_{re} , F_{re} denote the proposed ASRE module, temporal pooling, spatial pooling, FC layer and 2d convolution with filter size of 1×1 . $X_{re_ID} \in \mathbb{R}^{T \times \frac{C}{4} \times H \times W}$, $f_{re_ID} \in \mathbb{R}^{T \times \frac{C}{4}}$. The ID-relevant attribute feature and the corresponding prediction can be expressed as:

$$\begin{aligned}
 f_{re_attr} &= F_{TP}(F_{SP}(X_{re_attr})), \\
 p_{re} &= W_1(\sigma(BN(W_2(f_{re_attr})))),
 \end{aligned} \tag{3}$$

where W_1 , W_2 are FC layers, σ is the activation function and the BN is the batch norm layer. $f_{re_attr} \in \mathbb{R}^{\frac{C}{4}}$ is the ID-relevant attribute feature and $p_{re} \in \mathbb{R}^{D_{re_attr}}$ is the attribute prediction where D_{re_attr} is the number of ID-relevant attributes to predict. In this paper, we transform a multi-category classification problem into multiple binary-category classification problems. For example, there are five kinds of top color attribute on MARS dataset, we transform them into five binary classification problems. Therefore, D_{re_attr} equals to the sum of number of all categories for all attributes.

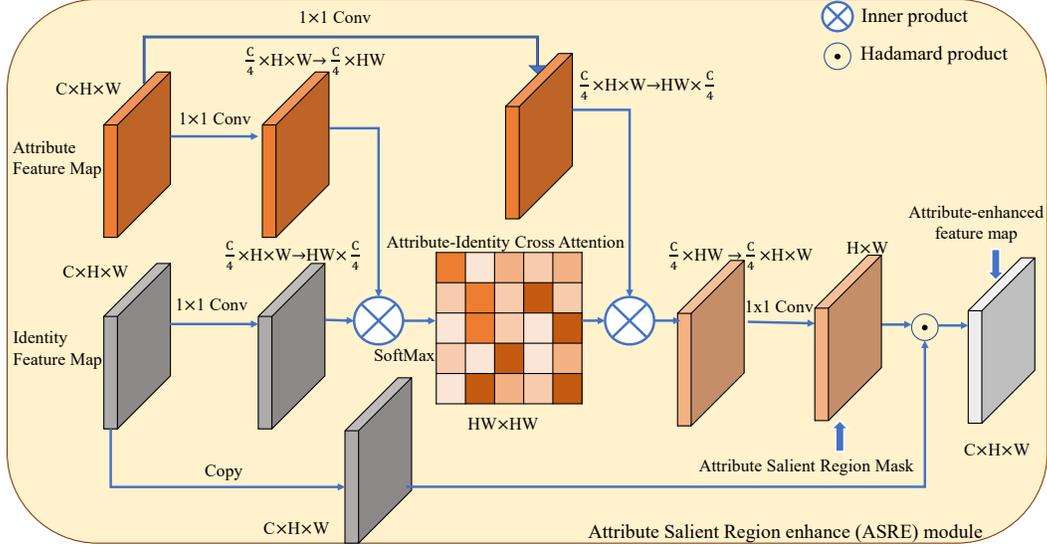


Fig. 4: Illustration of the Attribute Salient Region Enhance (ASRE) module. The response degree of each pixel in the identity feature map to each pixel in the attribute feature map can be calculated by cross attention mechanism. Then the obtained Attribute-Identity attention map multiplies with the projected attribute feature map to get the response value of the attribute on each pixel of the identity feature map. Finally, the final attribute salient region mask is calculated from response value map.

The framework of ID-irrelevant attribute branch is same to the ID-relevant attribute branch. The ID-irrelevant attribute branch can be formulated as following:

$$\begin{aligned} \mathbf{X}_{ir_ID} &= ASRE(\mathbf{X}_{ir_attr}, F_{ir}(\mathbf{X}_{ID})), \\ \mathbf{f}_{ir_ID} &= W_{ir}(F_{TP}(F_{SP}(\mathbf{X}_{ir_ID}))), \end{aligned} \quad (4)$$

where F_{ir} and W_{ir} are a 2d convolution layer with filter size of 1×1 and FC layer. $\mathbf{X}_{ir_ID} \in \mathbb{R}^{T \times \frac{C}{4} \times H \times W}$, $\mathbf{f}_{ir_ID} \in \mathbb{R}^{T \times \frac{C}{4}}$. The ID-irrelevant attribute prediction can be expressed as:

$$\begin{aligned} \mathbf{f}_{ir_attr} &= F_{TP}(F_{SP}(\mathbf{X}_{ir_attr})), \\ \mathbf{p}_{ir} &= W_3(\sigma(BN(W_4(\mathbf{f}_{ir_attr})))), \end{aligned} \quad (5)$$

where W_3 , W_4 are FC layers. $\mathbf{f}_{ir_attr} \in \mathbb{R}^{\frac{C}{4}}$ is the ID-irrelevant attribute feature and $\mathbf{p}_{ir} \in \mathbb{R}^{D_{ir_attr}}$ is the ID-irrelevant attribute prediction where D_{ir_attr} is the number of ID-irrelevant attributes to predict.

In order to keep the dimension consistent, we also reduce the dimension of original identity feature map \mathbf{X}_{ID} and calculate the corresponding feature \mathbf{f}_{ID} :

$$\begin{aligned} \mathbf{X}_{ID_ID} &= F_{ID}(\mathbf{X}_{ID}), \\ \mathbf{f}_{ID} &= F_{TP}(F_{SP}(\mathbf{X}_{ID_ID})), \end{aligned} \quad (6)$$

where F_{ID} is a 2d convolution layer with filter size of 1×1 , $\mathbf{X}_{ID_ID} \in \mathbb{R}^{T \times \frac{C}{4} \times H \times W}$ and $\mathbf{f}_{ID} \in \mathbb{R}^{T \times \frac{C}{4}}$.

Combining Eq. (2), (3), (4), (5) and (6), the final feature are calculated as:

$$\mathbf{f}_{final} = F_{fusion}(\mathbf{f}_{re_attr}, \mathbf{f}_{re_ID}, \mathbf{f}_{ID}, \mathbf{f}_{ir_attr}, \mathbf{f}_{ir_ID}), \quad (7)$$

where $\mathbf{f}_{final} \in \mathbb{R}^c$ and the F_{fusion} is the fusion module which will be explained in Section III-D.

B. Attribute Salient Region Enhance(ASRE) Module

Background clutter often interferes in person Re-ID task, so some work are devoted to using attention for enhancing the

features of human region [44], [45]. However, their attention is calculated with the input image itself. Inspired by [71], middle-level semantic labels of covariates often indicate the spatial locations of covariates. Therefore, additional middle-level semantic labels may indicate the spatial location of the person and make the attention more precise in person Re-ID. In this paper, While providing attribute information, the attribute annotation can also implicitly contain the relevant spatial cues of human body. Therefore the Attribute Salient Region Enhance (ASRE) module is proposed to find salient regions of pedestrian with attribute labels.

The procedure of ID-relevant attribute branch and ID-irrelevant attribute branch is similar, so we explain the ASRE module of ID-relevant attribute branch as an example.

Let's denote the $F_{re}(\mathbf{X}_{ID})$ as \mathbf{X}_I and \mathbf{X}_{re_attr} as \mathbf{X}_A which represent the identity feature map and the ID-relevant attribute feature map in Eq. (6), (3). Self-attention is a powerful long-range relationship capture mechanism [46], we modified it into a cross-attention mechanism to capture the relationship between two feature maps. As shown in Fig. 4, the Attribute-Identity cross attention map can be calculated as:

$$Att(\mathbf{X}_I, \mathbf{X}_A) = Softmax\left(\frac{\sigma(F_Q(\mathbf{X}_I))^T \sigma(F_K(\mathbf{X}_A))}{\sqrt{d_k}}\right), \quad (8)$$

where F_Q and F_K are combination of convolution layers with filter size of 1×1 and dimension merge operation. σ is the activation function and $\sqrt{d_k}$ is the scaling factor to prevent excessive values. In practice, σ is the ReLU function and d_k equals to the feature dimension after the 1×1 convolution. This attention map indicates the response degree of each pixel of the identity feature map to each pixel of the attribute feature map. Then the salient-region mask can be formulated as:

$$\begin{aligned} \mathbf{V} &= F_V(\mathbf{X}_A), \\ \mathbf{X}_V &= OP_{split}(Att(\mathbf{X}_I, \mathbf{X}_A)\mathbf{V}), \\ \mathbf{M}_{ASR} &= F_{mask}(BN(\mathbf{X}_V)), \end{aligned} \quad (9)$$

where F_V is combination of convolution layers with filter size of 1×1 and dimension merge operation, OP_{split} is the operation of changing the dimension from $\frac{C}{4} \times HW$ to $\frac{C}{4} \times H \times W$ and F_{mask} is a 2d convolution layer with filter size of 1×1 . $\mathbf{V} \in \mathbb{R}^{HW \times \frac{C}{4}}$, $\mathbf{X}_V \in \mathbb{R}^{\frac{C}{4} \times H \times W}$ and $\mathbf{M}_{ASR} \in \mathbb{R}^{H \times W}$

Combining with Eq. (8), (9), the attribute-enhanced feature map and final output of the ASRE module can be written as:

$$\begin{aligned} \mathbf{X}_{enhanced} &= \mathbf{M}_{ASR} \otimes \mathbf{X}_I, \\ \mathbf{X}_{re_ID} &= \mathbf{X}_I + \alpha * \mathbf{X}_{enhanced}, \end{aligned} \quad (10)$$

where \otimes is Hadamard product and α is a learnable parameter. The default value of α is 1.0 for both branches.

C. Pose & Motion-Invariant Triplet Loss (PMI Triplet Loss)

The common triplet loss aims at enlarging the distance between positive and negative pairs. However, it can not narrow the intra-class distance caused by the variation introduced by pose and motion. Some works [49], [45], [67] tried to solve this problem in image-based person Re-ID. In this work, with the obtained ID-irrelevant attribute feature \mathbf{f}_{ir_attr} , we can measure the difference of pose and motion between positive sample pairs.

Considering that there are K samples $\{I_1, I_2, \dots, I_K\}$ of one person in a mini-batch and the ID-irrelevant attribute features of them are $\{\mathbf{f}_{ir_attr}^1, \mathbf{f}_{ir_attr}^2, \dots, \mathbf{f}_{ir_attr}^K\}$ calculated in Eq. (5). For every anchor sample I_A , the triplet pairs are generated as following:

$$\begin{aligned} AP &= \arg \min_{1 \leq i \leq K, i \neq anchor} (\|\mathbf{f}_{ir_attr}^i - \mathbf{f}_{ir_attr}^A\|_2), \\ AN &= \arg \max_{1 \leq i \leq K, i \neq anchor} (\|\mathbf{f}_{ir_attr}^i - \mathbf{f}_{ir_attr}^A\|_2), \end{aligned} \quad (11)$$

and the triplet pair for anchor sample I_A is $\{I_A, I_{AP}, I_{AN}\}$. Supposing there N triplets in a mini-batch which are denoted as $\{T_i | T_i = (\mathbf{f}_{final}^{A_i}, \mathbf{f}_{final}^{AP_i}, \mathbf{f}_{final}^{AN_i})\}$ where \mathbf{f}_{final} is the final feature calculated in Eq. (7). Finally, the PMI triplet loss can be formulated as:

$$\mathcal{L}_{PMI} = \frac{1}{N} \sum_{i=1}^N \max(d_i^- - d_i^+, 0), \quad (12)$$

where N is the number of samples within a mini-batch and

$$\begin{aligned} d_i^- &= 1 - \frac{\mathbf{f}_{final}^{A_i} \mathbf{f}_{final}^{AN_i}}{|\mathbf{f}_{final}^{A_i}| |\mathbf{f}_{final}^{AN_i}|}, \\ d_i^+ &= 1 - \frac{\mathbf{f}_{final}^{A_i} \mathbf{f}_{final}^{AP_i}}{|\mathbf{f}_{final}^{A_i}| |\mathbf{f}_{final}^{AP_i}|}, \end{aligned} \quad (13)$$

denote the distance of the final feature between (I_A, I_{AN}) and (I_A, I_{AP}) .

D. Feature Fusion Strategy

As shown in Fig. 4, in the end of network pipeline, features from five branches need to be fused into one identity feature. The proposed two fusion strategies are shown in Fig. 5. For both strategies, \mathbf{f}_{re_ID} , \mathbf{f}_{ID} , \mathbf{f}_{ir_ID} need to go through a

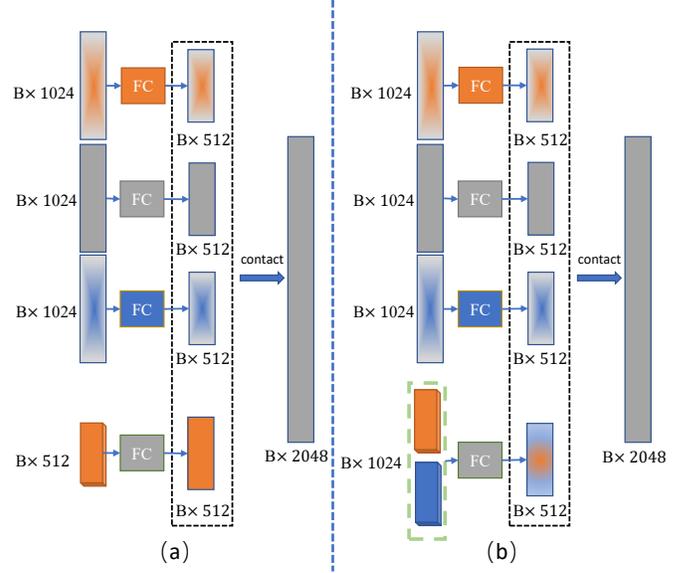


Fig. 5: Fusion strategy of features of five branches. Each strip corresponds to the same one of Fig. 3. The orange-grey mixed strip, grey strip, blue-grey mixed strip, orange strip and blue strip denote ID-relevant enhanced identity feature, original identity feature, ID-irrelevant enhanced identity feature, ID-relevant attribute feature and ID-irrelevant attribute feature respectively. All attribute enhanced features go through a Fully connect layer to reduce dimensions. In fusion strategy (a), only ID-relevant attribute feature are fused into final feature while in fusion strategy (b) we do a early fusion of ID-relevant attribute feature and ID-irrelevant attribute feature.

FC layer for dimension reduction. For strategy (a), only the feature of ID-relevant attributes \mathbf{f}_{re_ID} are fused in the final feature. For strategy (b), we contact features of ID-relevant and ID-irrelevant attributes into one attribute feature. Then the attribute goes through a FC layer and is contacted with other features.

Generally, the feature of ID-irrelevant is considered harmful to the Re-ID task. However, we argue that the ID-irrelevant features including pose and motion information can help to filter ID-relevant features. The experiments confirmed our supposition.

E. The Overall Loss

As illustrated in Fig. 4, our losses consist of five parts, i.e. weighted regularization triplet loss, center loss, cross entropy loss, binary cross entropy loss and the pose&motion-invariant loss.

Weighted regularization triplet loss was proposed in [50] which can be expressed as:

$$\mathcal{L}_{wrt} = \sum_{i=1}^N \log(1 + \exp(\sum_j (w_{ij}^p d_{ij}^p) - \sum_k (w_{ik}^n d_{ik}^n))) \quad (14)$$

$$w_{ij}^p = \frac{\exp(d_{ij}^p)}{\sum_{d_{ij}^p \in \mathcal{P}_i} \exp(d_{ij}^p)}, w_{ik}^n = \frac{\exp(-d_{ik}^n)}{\sum_{d_{ik}^n \in \mathcal{N}_i} \exp(-d_{ik}^n)}$$

where (i, j, k) represents a triplet within each training batch. For anchor i , \mathcal{P}_i is the corresponding positive set and \mathcal{N}_i is

the negative set. d_{ij}^p/d_{ik}^n represents the pairwise distance of a positive/negative sample pair.

Center loss was proposed in [51] which aims at minimizing the intra-class variations. It can be formulated as:

$$\mathcal{L}_{cent} = \frac{1}{2} \sum_{i=1}^N \|\mathbf{f}_i - \mathbf{c}_{y_i}\|_2, \quad (15)$$

where $\mathbf{c}_{y_i} \in \mathbb{R}^c$ denotes the center of y_i th person's features and is a learnable parameter.

Cross entropy loss has been used in Re-ID task widely. In this paper we used the cross entropy loss with label smooth as [52]. Supposing there are M persons in the training dataset, for the obtained final feature \mathbf{f}_{final} in Eq.(7), we calculated the prediction probability of sample i as:

$$\mathbf{p}_i = W_{ID}(BN(\mathbf{f}_{final,i})), \quad (16)$$

where BN is a batchnorm layer, W_{ID} is a FC layer and $\mathbf{p}_i \in \mathbb{R}^M$. Then the cross entropy loss with label smooth can be expressed as:

$$\mathcal{L}_{xent} = \sum_{i=1}^N \sum_{j=1}^M -q_{ij} \log(\mathbf{p}_{ij}) \begin{cases} q_{ij} = 1 - \frac{N-1}{N} \epsilon, y_i = j \\ q_{ij} = \frac{\epsilon}{N}, otherwise \end{cases}, \quad (17)$$

where y_i is the ID label of sample i .

The binary cross entropy loss is used for attribute prediction. We merge the obtained attribute prediction probability \mathbf{p}_{re} and \mathbf{p}_{ir} in Eq. (3), (5) and denote them as \mathbf{p}_{attr} . The i th prediction of sample i is denoted as $\mathbf{p}_i^{attr} \in \mathbb{R}^{D_{attr}}$ where $D_{attr} = D_{re_attr} + D_{ir_attr}$. Then the binary cross entropy loss of attributes are formulated as:

$$\mathcal{L}_{BCE} = \sum_{i=1}^N \sum_{j=1}^{D_{attr}} -y_{ij}^{attr} \log(\mathbf{p}_{ij}^{attr}), \quad (18)$$

where $y_{ij}^{attr} \in \{0,1\}$ is the one-hot label of j th attribute of i th sample.

Finally, combining Eq. (12), (14), (15), (17), (18), the loss function of the model can be written as:

$$\mathcal{L} = \mathcal{L}_{xent} + \mathcal{L}_{wrt} + \lambda_{cent} \mathcal{L}_{cent} + \lambda_{BCE} \mathcal{L}_{BCE} + \lambda_{PMI} \mathcal{L}_{PMI}, \quad (19)$$

where λ_{cent} , λ_{BCE} , λ_{PMI} are hyper-parameters.

IV. EXPERIMENT

A. Datasets and Evaluation Protocols

We evaluate the proposed model on two large-scale video person Re-ID datasets: DukeMTMC-VideoReID [42] and MARS [41]. DukeMTMC-VideoReID comprises around 4,832 videos from 1,812 identities and we use *Duke* for abbreviation. There are 702 and 702 identities for training and testing in Duke as well as 408 identities for distraction. The bounding boxes are annotated manually. MARS is the largest dataset for video-based person Re-ID and has 17,503 sequences of 1,261 identities and 3,248 dis-tractor sequences. Among 1,261 identities, 625 identities are for training and the other 636 identities are for testing. The annotation of attributes we use is provided by Chen et al. [47].

TABLE I: Performance(%) comparison with related works on MARS. ‘‘a’’ and ‘‘b’’ denote the fusion strategy. PMI means the model is trained with PMI triplet loss. ‘‘Mixing’’ means the result is tested under the situation of mixing the query set and the gallery set as the new gallery set. AP3D and TCLNet are tested under this special situation according to their code to avoid the lack of corresponding gallery samples. The highest and second highest performance are **bold** and underlined under usual setup. The highest performance are **bold** under Mixing setup.

Setup	Method	mAP	Rank-1	Rank-5	Rank-20
Usual	BoW+kissme [41]	15.5	30.6	46.2	59.2
	IDE+XQDA [41]	47.6	65.3	82.0	89.0
	Wu et al. [87]	-	73.5	85.0	97.5
	MG-TCN [73]	77.7	87.0	95.1	98.2
	STA [59]	80.8	86.3	95.7	98.1
	Attribute [34]	78.2	87.0	95.8	<u>98.7</u>
	HMN [72]	82.6	88.5	96.2	98.1
	VRSTC [60]	82.3	88.5	96.5	97.4
	GLTR [61]	78.5	87.0	95.8	98.2
	COSAM [44]	79.9	84.9	95.5	97.9
	AMEM [35]	79.3	86.7	94.0	97.1
	MGRA [43]	85.9	88.8	<u>97.0</u>	98.5
	STGCN [36]	83.7	89.9	96.4	98.3
	FA [38]	82.9	90.2	96.6	-
	MGH [40]	85.8	90.0	96.7	98.5
	ASTA [82]	84.1	90.4	97.0	98.8
	TALNet [65]	82.3	89.1	96.1	98.5
	PhD [85]	86.2	88.9	<u>97.0</u>	98.6
	ASANet-b(Ours)	86.6	<u>90.3</u>	96.8	98.4
	ASANet-a-PMI(Ours)	<u>86.3</u>	89.5	97.1	<u>98.7</u>
Mixing	TCLNet [37]	85.1	89.8	-	-
	AP3D [39]	85.6	90.7	-	-
	BiCnet [83]	86.0	90.2	-	-
	ASANet-b(Ours)	86.0	91.1	97.0	98.4
	ASANet-a-PMI(Ours)	86.0	90.6	97.1	98.7

Cumulative Matching Characteristic (CMC) is widely used to evaluate the performance of person Re-ID models and we adopt the CMC to evaluate our model. Also, the mean Average Precision (mAP) is also adopted for evaluation.

B. Implementation Details

The baseline of our net is AGW-Net [50] which is based on ResNet-50 [11] pre-trained on ImageNet [54]. The backbone part in Fig. 3 only contains the first four layers of ResNet-50. As the same as [52], each image is resized to 256×128 , random horizontal flip, random crop and random erasing are adopted as data augmentation. We applied the constrain random sampling strategy [55] to randomly sample $T = 6$ frames from every clip. We train our network for 700 epochs in total, with the initial learning rate of 0.0003 and decayed it at the 100, 250, 350 epoch. For parameters of the ASA-network, we use Adam [57] algorithm with weight decay of $5e-4$ and for parameters of the center loss in Equation (15), we use SGD [56] algorithm with learning rate of 0.5. For every mini-batch, we sample eight identities, each with four tracklets, to form a mini-batch of size $8 \times 4 \times 6 = 192$ images. λ_{cent} , λ_{BCE} are set to 1.5 and 0.0005. λ_{PMI} is set to 0.005 initially and turn to 0.01 when \mathcal{L}_{BCE} is smaller than 0.15 which means the prediction of pose and motion is accurate enough. Our framework is implemented with Pytorch toolbox and two NVIDIA GTX 1080Ti GPUs.

TABLE II: Performance(%) comparison with related works on Duke. 'b' denote the fusion strategy. 'PMI' means the model is trained with PMI triplet loss.

Method	mAP	Rank-1	Rank-5	Rank-20
EUG [63]	78.3	83.6	94.6	97.6
ETAP-Net [63]	78.3	83.6	94.6	97.6
STE-NVAN [64]	93.5	95.2	-	-
VRSTC [60]	93.5	95.0	99.1	-
STA [59]	94.9	96.2	99.3	99.6
Wu et al. [62]	94.2	96.7	99.2	99.7
GLTR [61]	93.74	96.29	99.3	99.7
HMN [72]	95.1	96.3	99.2	99.8
STGCN [36]	95.7	97.3	99.3	99.7
AP3D [39]	96.1	97.2	-	-
TCLNet [37]	96.2	96.9	-	-
ASANet-b(Ours)	97.1	97.3	99.9	100.0
ASANet-b-PMI(Ours)	97.1	97.6	99.9	100.0

TABLE III: Cross-dataset evaluations between MARS and DukeMTMC-Video (Duke for short)(%).

Method	MARS->Duke			Duke->MARS		
	mAP	Rank-1	Rank-5	mAP	Rank-1	Rank-5
QAN [38]	33.0	38.9	59.3	24.1	43.3	58.8
AFA [38]	34.6	41.5	59.1	24.5	44.2	58.8
Ours	57.8	57.8	77.9	35.1	58.0	69.7

C. Comparison with the State-of-the-art Approaches

In this section, we compare the proposed method with state-of-the-art methods on two video person Re-ID benchmarks.

The results on MARS dataset are reported in Table I. There are 140 samples in the query set lack the corresponding samples in gallery set on MARS dataset. The usual testing setup is to remove these 140 samples from the query set. However, we find that Gu et al. [39] and Hou et al. [37] merge the query set and gallery set to form the new gallery set to avoid this problem according to their codes which may make the rank-1 value higher and the mAP value lower. We compare the results under two testing setups separately for fair. Some methods don't make their code available nor tell the detail in their papers, we make a reasonable guess that they tested their methods under the usual setup. Our experiments show that different fusion strategies and whether to use PMI triplet loss lead to different result on CMC. So we list the two combinations with the highest mAP here. The more analysis of fusion strategies and PMI triplet loss will be presented in the subsection of ablation study.

Under the normal setup, the mAP and Rank-1 accuracy of our method with fusion strategy (b) and without PMI triplet loss (ASANet-b) is 86.6% and 90.3% which have obvious improvement compared with SOTA methods. Especially, compared with the first work of using attribute [34], ASANet-b improves 8.4% on mAP and 3.3% on Rank-1. AMEM [35] considers both ID-relevant attributes and ID-irrelevant attributes as we do, we outperform AMEM [35] by 7.3% on mAP and 3.6% on Rank-1. The Rank-1 value of our ASANet-b is close to ASTA [82] which also put forward the 'motion' concept, but the mAP value of our method is 2.5% higher than it. To the same as our method, TALNet [65] uses the more accurate annotated attribute label and our ASANet-b is also 4.3% and 1.2% ahead of it on mAP and Rank-1. The mAP and Rank-1 value of our method with fusion

TABLE IV: Ablation Study of different branches on MARS(%).

Branches			Metrics			
Identity	ID-re.	ID-ir.	mAP	rank-1	rank-5	rank-20
✓			83.9	89.2	96.8	98.2
	✓		85.1	89.8	96.9	98.5
		✓	84.8	89.8	96.8	98.4
✓	✓		85.3	90.3	96.6	98.2
✓		✓	84.6	89.5	96.5	98.3
	✓	✓	85.8	89.9	96.8	98.3
✓	✓	✓	85.8	90.3	96.9	98.5

TABLE V: Ablation Study of fusion strategies, ASRE module and PMI triplet loss on MARS(%).

Fusion strategy	ASRE	PMI	mAP	rank-1	rank-5	rank-20
a			85.1	89.9	96.5	98.5
a		✓	85.5	90.2	97.0	98.5
a	✓		85.3	90.4	97.1	98.5
a	✓	✓	86.0	90.6	97.1	98.7
b			85.4	90.4	96.9	98.5
b		✓	85.7	90.5	97.3	98.5
b	✓		86.0	91.1	97.0	98.4
b	✓	✓	85.8	90.8	97.0	98.6

TABLE VI: Ablation Study of fusion strategies and PMI triplet loss on Duke(%).

Fusion strategy	PMI	mAP	rank-1	rank-5	rank-20
a		96.7	97.0	99.7	100.0
a	✓	97.1	97.4	99.9	100.0
b		97.1	97.3	99.9	100.0
b	✓	97.1	97.6	99.9	100.0

strategy (a) and with PMI triplet loss(ASANet-a-PMI) is a little lower than ASANet-b. But ASANet-a-PMI achieves the highest rank-5 and rank-20 score which is more important in practical application.

Under the mixing setup, the mAP value of both ASANet-b and ASANet-a-PMI is higher than other SOTA methods and the Rank-1 of ASANet-b also outperforms other methods.

The results on Duke dataset are reported in Table II. Our methods outperform the SOTA method TCLNet [37] by 0.9% on mAP which is a significant improvement since the mAP value is pretty high. It's worth noting that our method achieve 100% on rank-20 which indicates that our method can achieve 100% success in practical video-based person re-ID task (For example, tracking of escaped criminal) with manual screening.

D. Cross-dataset Evaluation

In real surveillance systems, due to the limitation of data annotation, models should be able to cope with data from different sources. Especially, our method has achieved excellent performance on Duke dataset. In order to verify that our method is not over fitting on Duke, we do the cross-dataset experiments. The results are shown in Table III. We use ASANet-b and ASANet-a-PMI, which have the best performance on MARS and Duke themselves to do the cross-dataset experiments. We also compare the results with the methods proposed by Chen et al. [38] and our methods show obvious advantages.



Fig. 6: Top 1 retrieval results of the ASANet-b and ASANet-b-PMI. In exemplar sequence (a), the retrieved sequence without PMI triplet loss and the query sequence have almost the same pose and action, and their clothes and background are also extremely similar. But it can be seen from the head that they are not the same person. Although the sequences retrieved with PMI triplet loss are different in pose and background, they are the same person. Sequence (b) shows that our model is also robust to pose change caused by camera view difference.

TABLE VII: Ablation study of attribute supervision(%).

Method	mAP	Rank-1	Rank-5	Rank-20
ASANet-b w/o BCE loss	85.5	89.8	97.0	98.6
ASANet-b with BCE loss	86.0	91.1	97.0	98.4

E. Ablation Study

In this subsection, we conduct experiments to verify the effectiveness of the proposed methods. Experiments are mainly conducted on MARS dataset since it has the largest amount of sequences. Also, our experiments are done under the mixing setup.

1) *Effects of each branch*: The result in Table IV shows the influence of each branch. ID-re. and ID-ir. denote the ID-relevant attribute enhanced feature and the ID-irrelevant attribute enhanced feature respectively. Note that we only use the enhanced feature and don't fuse attribute features here. Also, all the experiments in Table IV are done with ASRE module and PMI triplet loss. All of the branches make sense. It can be seen that the Identity branch and ID-re branch plays a more important role in improving the performance which is congenial with common sense.

2) *Effects of the fusion strategies, ASRE module and the PMI triplet loss*: On the basis of the three-branch framework, we explore the role of the fusion strategies, the ASRE module and the PMI triplet loss. The results are reported in Table V. It can be seen that whether under the fusion strategy (a) or under the fusion strategy (b), ASRE module or PMI triplet loss can boost the performance. Under the fusion strategy

a, the joint use of ASRE module and PMI triplet loss can further improve the performance. On the contrary, the joint use of ASRE module and PMI triplet loss will decrease the performance of the model under the fusion strategy (b). In order to explore the reason of this phenomenon, we observe the bad cases which is identified successfully with adding ASRE module only and failed with adding both ASRE and PMI. All of these bad cases retrieved the dis-tractor sequences which contain no person. In the training set, we have no dis-tractor sequences and the ASANet-b-PMI maybe is too complex to generalize to dis-tractor sequences which have no person.

To further prove the effectiveness of our proposed PMI triplet loss, we do the ablation study of PMI triplet loss on the Duke dataset which doesn't have sequences without pedestrian. The results are shown in Table VI.

3) *Effects of the Attribute Assistance*: In order to explore whether attribute plays a role or whether our framework itself plays a role, we do the ablation experiments on the BCE loss. The results are shown in Table VII. It can be seen that the attribute supervision can make significant improvement of mAP and Rank-1 value.

F. Visualization

In order to illustrate the effectiveness of our method, we do some visualization.

Firstly, we visualize the effect of PMI triplet loss. According to the results on Duke dataset shown in Table VI, the improvement caused by PMI triplet loss is mainly reflected in

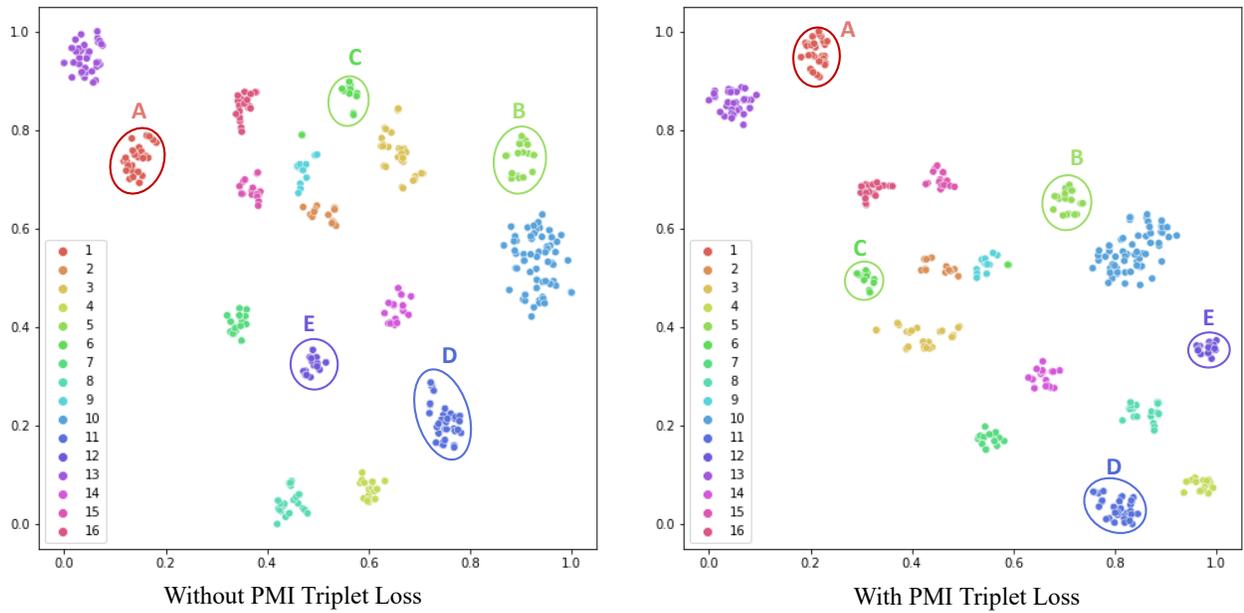


Fig. 7: Visualization of feature distribution of the first 16 subjects on MARS gallery dataset using t-SNE. Subjects whose numbers of samples less than 10 is skipped because they are unsuitable for observation. It can be seen that the proposed PMI triplet loss can narrow the intra-class distance to a certain extent. For example, features of subject A, D, E can be divided into more sub-classes without PMI triplet loss than with PMI triplet loss.

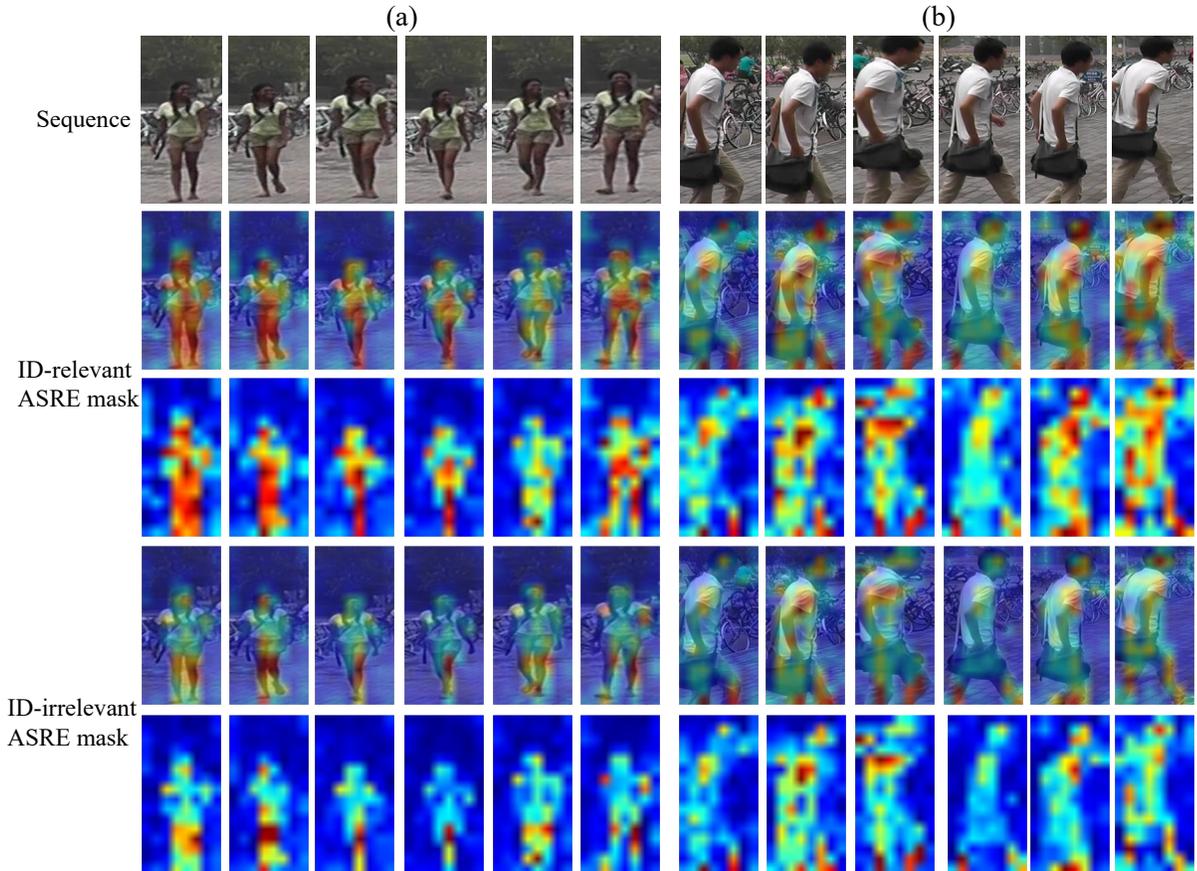


Fig. 8: Attention mask generated by ASRE module. It can be seen that the mask regions generated by the two branches are roughly the same and can accurately cover the human body regions. But they differs in details. ID-relevant ASRE mask pay more attention to the detail of all the human body while ID-irrelevant ASRE mask only pay more attention to legs and joints of human body. ID-irrelevant ASRE mask can outline the edge of human body and reduce the influence of background more accurately.

rank-1 value. Therefore, we show some top-1 retrieval results in Fig. 6. It can be seen that PMI triplet loss can cope with the change of pose and motion perfectly. Especially in Fig. 6 (a), the retrieved sequence share almost the same pose and motion with the query sequence. Also, we choose the first 16 subjects (ordered by ID number) on MARS gallery dataset to visualize the feature distribution by t-SNE [66]. Subjects whose numbers of samples less than 10 is skipped because they are unsuitable for observation. The results are shown in Fig. 7. It can be seen that PMI triplet loss can effectively narrow the intra-class distance.

Secondly, to explore the efficacy of proposed ASRE module, we visualized the mask generated by the ASRE module. As illustrated in Fig. 8, the generated mask can mainly focus on the pedestrian body region and pay more attention to the area which has salient details. It's consistent with our design that mask generated by ID-relevant branch pay more attention to the region of human ID-relevant attributes such as the shoulder bag in Fig. 8 (b) and mask generated by ID-irrelevant branch can outline the edge of human body more accurately and reduce the influence of background such as Fig. 8 (a).

V. CONCLUSION

In this work, we propose a comprehensive study on attribute-assisted video person re-identification. We introduced both ID-relevant attributes and ID-irrelevant attributes. To integrate attribute information into the Re-ID task we propose the ASA-Net and PMI triplet loss. The former predicts attributes, extract identity features at the same time and further enhance them by the attribute-salient regions learned from the ASRE module. With the extracted ID-irrelevant attribute feature, we form the PMI triplet loss to narrow the intra-class distance caused by the change of pose and motion. We conducted experiments on two challenging datasets and outperforms the state-of-the-art methods. Some visualizations are also presented to illustrate the effectiveness of our method.

REFERENCES

- [1] Chunxiao Liu, Shaogang Gong, Chen Change Loy, and Xinggang Lin, "Person re-identification: What features are important?," in *ECCV*. Springer, 2012, pp. 391–401.
- [2] Igor Kviatkovsky, Amit Adam, and Ehud Rivlin, "Color invariants for person reidentification," *IEEE TPAMI*, vol. 35, no. 7, pp. 1622–1634, 2012.
- [3] Liming Zhao, Xi Li, Yueting Zhuang, and Jingdong Wang, "Deeply-learned part-aligned representations for person re-identification," in *ICCV*, 2017, pp. 3219–3228.
- [4] Martin Hirzer, Csaba Beleznai, Peter M Roth, and Horst Bischof, "Person re-identification by descriptive and discriminative classification," in *SCIA*. Springer, 2011, pp. 91–102.
- [5] Song Bai, Xiang Bai, and Qi Tian, "Scalable person re-identification on supervised smoothed manifold," in *CVPR*, 2017, pp. 2530–2539.
- [6] Fei Xiong, Mengran Gou, Octavia Camps, and Mario Szaier, "Person re-identification using kernel-based metric learning methods," in *ECCV*. Springer, 2014, pp. 1–16.
- [7] Shengcai Liao and Stan Z Li, "Efficient psd constrained asymmetric metric learning for person re-identification," in *ICCV*, 2015, pp. 3685–3693.
- [8] Dapeng Chen, Zejian Yuan, Badong Chen, and Nanning Zheng, "Similarity learning with spatial constraints for person re-identification," in *CVPR*, 2016, pp. 1268–1277.
- [9] Wei Li and Xiaogang Wang, "Locally aligned feature transforms across views," in *CVPR*, 2013, pp. 3594–3601.
- [10] Tetsu Matsukawa, Takahiro Okabe, Einoshin Suzuki, and Yoichi Sato, "Hierarchical gaussian descriptor for person re-identification," in *CVPR*, 2016, pp. 1363–1372.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [12] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *CVPR*, 2014, pp. 152–159.
- [13] Ejaz Ahmed, Michael Jones, and Tim K Marks, "An improved deep learning architecture for person re-identification," in *CVPR*, 2015, pp. 3908–3916.
- [14] Rahul Rama Varior, Mrinal Haloi, and Gang Wang, "Gated siamese convolutional neural network architecture for human re-identification," in *ECCV*. Springer, 2016, pp. 791–808.
- [15] Ji Lin, Liangliang Ren, Jiwen Lu, Jianjiang Feng, and Jie Zhou, "Consistent-aware deep learning for person re-identification in a camera network," in *CVPR*, 2017, pp. 5771–5780.
- [16] Taiqing Wang, Shaogang Gong, Xiatian Zhu, and Shengjin Wang, "Person re-identification by video ranking," in *ECCV*. Springer, 2014, pp. 688–703.
- [17] Dapeng Tao, Lianwen Jin, Yongfei Wang, Yuan Yuan, and Xuelong Li, "Person re-identification by regularized smoothing kiss metric learning," *IEEE TCSVT*, vol. 23, no. 10, pp. 1675–1685, 2013.
- [18] Niall McLaughlin, Jesus Martinez Del Rincon, and Paul Miller, "Re-current convolutional network for video-based person re-identification," in *CVPR*, 2016, pp. 1325–1334.
- [19] Zhen Zhou, Yan Huang, Wei Wang, Liang Wang, and Tieniu Tan, "See the forest for the trees: Joint spatial and temporal recurrent neural networks for video-based person re-identification," in *CVPR*, 2017, pp. 4747–4756.
- [20] Shuangjie Xu, Yu Cheng, Kang Gu, Yang Yang, Shiyu Chang, and Pan Zhou, "Jointly attentive spatial-temporal pooling networks for video-based person re-identification," in *ICCV*, 2017, pp. 4733–4742.
- [21] Yiheng Liu, Zhenxun Yuan, Wengang Zhou, and Houqiang Li, "Spatial and temporal mutual promotion for video-based person re-identification," in *AAAI*, 2019, vol. 33, pp. 8786–8793.
- [22] Yubin Deng, Ping Luo, Chen Change Loy, and Xiaoou Tang, "Pedestrian attribute recognition at far distance," in *ACM MM*, 2014, pp. 789–792.
- [23] Xihui Liu, Haiyu Zhao, Maoqing Tian, Lu Sheng, Jing Shao, Shuai Yi, Junjie Yan, and Xiaogang Wang, "Hydraplus-net: Attentive deep features for pedestrian analysis," in *ICCV*, 2017, pp. 350–359.
- [24] Xin Zhao, Liufang Sang, Guiguang Ding, Yuchen Guo, and Xiaoming Jin, "Grouping attribute recognition for pedestrian with joint recurrent learning," in *IJCAI*, 2018, pp. 3177–3183.
- [25] Lian Gao, Di Huang, Yuanfang Guo, and Yunhong Wang, "Pedestrian attribute recognition via hierarchical multi-task learning and relationship attention," in *ACM MM*, 2019, pp. 1340–1348.
- [26] Ryan Layne, Timothy M Hospedales, Shaogang Gong, and Q Mary, "Person re-identification by attributes," in *BMVC*, 2012, vol. 2, p. 8.
- [27] Annan Li, Luoqi Liu, Kang Wang, Si Liu, and Shuicheng Yan, "Clothing attributes assisted person reidentification," *IEEE TCSVT*, vol. 25, no. 5, pp. 869–878, 2014.
- [28] Chi Su, Fan Yang, Shiliang Zhang, Qi Tian, Larry Steven Davis, and Wen Gao, "Multi-task learning with low rank attribute embedding for multi-camera person re-identification," *IEEE TPAMI*, vol. 40, no. 5, pp. 1167–1181, 2017.
- [29] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li, "Transferable joint attribute-identity deep learning for unsupervised person re-identification," in *CVPR*, 2018, pp. 2275–2284.
- [30] Yutian Lin, Liang Zheng, Zhedong Zheng, Yu Wu, Zhilan Hu, Chenggang Yan, and Yi Yang, "Improving person re-identification by attribute and identity learning," *Pattern Recognition*, vol. 95, pp. 151–161, 2019.
- [31] Hefei Ling, Ziyang Wang, Ping Li, Yuxuan Shi, Jiazong Chen, and Fuhao Zou, "Improving person re-identification by multi-task learning," *Neurocomputing*, vol. 347, pp. 109–118, 2019.
- [32] Kai Han, Jianyuan Guo, Chao Zhang, and Mingjian Zhu, "Attribute-aware attention model for fine-grained representation learning," in *ACM MM*, 2018, pp. 2040–2048.
- [33] Wanru Song, Jieying Zheng, Yahong Wu, Changhong Chen, and Feng Liu, "A two-stage attribute-constraint network for video-based person re-identification," *IEEE Access*, vol. 7, pp. 8508–8518, 2019.
- [34] Yiru Zhao, Xu Shen, Zhongming Jin, Hongtao Lu, and Xian-sheng Hua, "Attribute-driven feature disentangling and temporal aggregation for video person re-identification," in *CVPR*, 2019, pp. 4913–4922.
- [35] Shuzhao Li, Huimin Yu, and Haoji Hu, "Appearance and motion enhancement for video-based person re-identification," in *AAAI*, 2020, vol. 34, pp. 11394–11401.

- [36] Jinrui Yang, Wei-Shi Zheng, Qize Yang, Ying-Cong Chen, and Qi Tian, "Spatial-temporal graph convolutional network for video-based person re-identification," in *CVPR*, 2020, pp. 3289–3299.
- [37] Ruibing Hou, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen, "Temporal complementary learning for video person re-identification," in *ECCV*. Springer, 2020, pp. 388–405.
- [38] Guangyi Chen, Yongming Rao, Jiwen Lu, and Jie Zhou, "Temporal coherence or temporal motion: Which is more critical for video-based person re-identification?," in *ECCV*. Springer, 2020, pp. 660–676.
- [39] Xinqian Gu, Hong Chang, Bingpeng Ma, Hongkai Zhang, and Xilin Chen, "Appearance-preserving 3d convolution for video-based person re-identification," in *ECCV*. Springer, 2020, pp. 228–243.
- [40] Yichao Yan, Jie Qin, Jiabin Chen, Li Liu, Fan Zhu, Ying Tai, and Ling Shao, "Learning multi-granular hypergraphs for video-based person re-identification," in *CVPR*, 2020, pp. 2899–2908.
- [41] Liang Zheng, Zhi Bie, Yifan Sun, Jingdong Wang, Chi Su, Shengjin Wang, and Qi Tian, "Mars: A video benchmark for large-scale person re-identification," in *ECCV*. Springer, 2016, pp. 868–884.
- [42] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," in *ECCV*. Springer, 2016, pp. 17–35.
- [43] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen, "Multi-granularity reference-aided attentive feature aggregation for video-based person re-identification," in *CVPR*, 2020, pp. 10407–10416.
- [44] Arulkumar Subramaniam, Athira Nambiar, and Anurag Mittal, "Co-segmentation inspired attention networks for video-based person re-identification," in *ICCV*, 2019, pp. 562–572.
- [45] Jiabin Chen, Jie Qin, Yichao Yan, Lei Huang, Li Liu, Fan Zhu, and Ling Shao, "Deep local binary coding for person re-identification by delving into the details," in *ACM MM*, 2020, pp. 3034–3043.
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [47] Zhiyuan Chen, Annan Li, and Yunhong Wang, "A temporal attentive approach for video-based pedestrian attribute recognition," in *PRCV*. Springer, 2019, pp. 209–220.
- [48] Pu Chen, Xinyi Xu, and Cheng Deng, "Deep view-aware metric learning for person re-identification," in *IJCAI*, 2018, pp. 620–626.
- [49] Zhihui Zhu, Xinyang Jiang, Feng Zheng, Xiaowei Guo, Feiyue Huang, Xing Sun, and Weishi Zheng, "Aware loss with angular regularization for person re-identification," in *AAAI*, 2020, vol. 34, pp. 13114–13121.
- [50] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi, "Deep learning for person re-identification: A survey and outlook," *IEEE TPAMI*, 2021.
- [51] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao, "A discriminative feature learning approach for deep face recognition," in *ECCV*. Springer, 2016, pp. 499–515.
- [52] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang, "Bag of tricks and a strong baseline for deep person re-identification," in *CVPR*, 2019, pp. 0–0.
- [53] Jianfu Zhang, Li Niu, and Liqing Zhang, "Person re-identification with reinforced attribute attention selection," *IEEE TIP*, vol. 30, pp. 603–616, 2020.
- [54] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*. Ieee, 2009, pp. 248–255.
- [55] Shuang Li, Slawomir Bak, Peter Carr, and Xiaogang Wang, "Diversity regularized spatiotemporal attention for video-based person re-identification," in *CVPR*, 2018, pp. 369–378.
- [56] Léon Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*, pp. 177–186. Springer, 2010.
- [57] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [58] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang, "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," in *ECCV*, 2018, pp. 480–496.
- [59] Yang Fu, Xiaoyang Wang, Yunchao Wei, and Thomas Huang, "Sta: Spatial-temporal attention for large-scale video-based person re-identification," in *AAAI*, 2019, vol. 33, pp. 8287–8294.
- [60] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen, "Vrstc: Occlusion-free video person re-identification," in *CVPR*, 2019, pp. 7183–7192.
- [61] Jianing Li, Jingdong Wang, Qi Tian, Wen Gao, and Shiliang Zhang, "Global-local temporal representations for video person re-identification," in *ICCV*, 2019, pp. 3958–3967.
- [62] Yiming Wu, Omar El Farouk Bourahla, Xi Li, Fei Wu, Qi Tian, and Xue Zhou, "Adaptive graph representation learning for video person re-identification," *IEEE TIP*, vol. 29, pp. 8821–8830, 2020.
- [63] Yu Wu, Yutian Lin, Xuanyi Dong, Yan Yan, Wanli Ouyang, and Yi Yang, "Exploit the unknown gradually: One-shot video-based person re-identification by stepwise learning," in *CVPR*, 2018, pp. 5177–5186.
- [64] Chih-Ting Liu, Chih-Wei Wu, Yu-Chiang Frank Wang, and Shao-Yi Chien, "Spatially and temporally efficient non-local attention network for video-based person re-identification," *arXiv preprint arXiv:1908.01683*, 2019.
- [65] Jiawei Liu, Xierong Zhu, and Zheng-Jun Zha, "Temporal attribute-appearance learning network for video-based person re-identification," *arXiv preprint arXiv:2009.04181*, 2020.
- [66] Laurens Van der Maaten and Geoffrey Hinton, "Visualizing data using t-sne," *JMLR*, vol. 9, no. 11, 2008.
- [67] Liang Zheng, Yujia Huang, Huchuan Lu, and Yi Yang, "Pose-invariant embedding for deep person re-identification," *IEEE TIP*, vol. 28, no. 9, pp. 4500–4509, 2019.
- [68] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian, "Person re-identification in the wild," in *CVPR*, 2017, pp. 1367–1376.
- [69] Alexander Hermans, Lucas Beyer, and Bastian Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [70] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang, "Beyond triplet loss: a deep quadruplet network for person re-identification," in *CVPR*, 2017, pp. 403–412.
- [71] Tianrui Chai, Xinyu Mei, Annan Li, and Yunhong Wang, "Semantically-guided disentangled representation for robust gait recognition," in *ICME*. IEEE, 2021, pp. 1–6.
- [72] Zhikang Wang, Lihuo He, Xiaoguang Tu, Jian Zhao, Xinbo Gao, Shengmei Shen, and Jiashi Feng, "Robust video-based person re-identification by hierarchical mining," *IEEE TCSVT*, 2021.
- [73] Peike Li, Pingbo Pan, Ping Liu, Mingliang Xu, and Yi Yang, "Hierarchical temporal modeling with mutual distance matching for video based person re-identification," *IEEE TCSVT*, vol. 31, no. 2, pp. 503–511, 2020.
- [74] Wei Zhang, Shengnan Hu, Kan Liu, and Zhengjun Zha, "Learning compact appearance representation for video-based person re-identification," *IEEE TCSVT*, vol. 29, no. 8, pp. 2442–2452, 2018.
- [75] Zheng Liu, Yunhong Wang, and Annan Li, "Hierarchical integration of rich features for video-based person re-identification," *IEEE TCSVT*, vol. 29, no. 12, pp. 3646–3659, 2018.
- [76] Yuxuan Shi, Zhen Wei, Hefei Ling, Ziyang Wang, Jialie Shen, and Ping Li, "Person retrieval in surveillance videos via deep attribute mining and reasoning," *IEEE TMM*, 2020.
- [77] Chiat-Pin Tay, Sharmili Roy, and Kim-Hui Yap, "Aanet: Attribute attention network for person re-identifications," in *CVPR*, 2019, pp. 7134–7143.
- [78] Huafeng Li, Shuanglin Yan, Zhengtao Yu, and Dapeng Tao, "Attribute-identity embedding and self-supervised learning for scalable person re-identification," *IEEE TCSVT*, vol. 30, no. 10, pp. 3472–3485, 2019.
- [79] Zheng Wang, Junjun Jiang, Yang Wu, Mang Ye, Xiang Bai, and Shin'ichi Satoh, "Learning sparse and identity-preserved hidden attributes for person re-identification," *IEEE TIP*, vol. 29, pp. 2013–2025, 2019.
- [80] Huafeng Li, Zhenyu Kuang, Zhengtao Yu, and Jiebo Luo, "Structure alignment of attributes and visual features for cross-dataset person re-identification," *Pattern Recognition*, vol. 106, pp. 107414, 2020.
- [81] Cheng Yan, Guansong Pang, Xiao Bai, Changhong Liu, Ning Xin, Lin Gu, and Jun Zhou, "Beyond triplet loss: person re-identification with fine-grained difference-aware pairwise loss," *IEEE TMM*, 2021.
- [82] Xierong Zhu, Jiawei Liu, Haoze Wu, Meng Wang, and Zheng-Jun Zha, "Asta-net: Adaptive spatio-temporal attention network for person re-identification in videos," in *ACM MM*, 2020, pp. 1706–1715.
- [83] Ruibing Hou, Hong Chang, Bingpeng Ma, Rui Huang, and Shiguang Shan, "Bicnet-tks: Learning efficient spatial-temporal representation for video person re-identification," in *CVPR*, 2021, pp. 2014–2023.
- [84] Zimo Liu, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang, "Person reidentification by joint local distance metric and feature transformation," *IEEE TNNLS*, vol. 30, no. 10, pp. 2999–3009, 2019.
- [85] Jianan Zhao, Fengliang Qi, Guangyu Ren, and Lin Xu, "Phd learning: Learning with pompeiu-hausdorff distances for video-based vehicle re-identification," in *CVPR*, June 2021, pp. 2225–2235.
- [86] Xun Gong, Zu Yao, Xin Li, Yueqiao Fan, Bin Luo, Jianfeng Fan, and Boji Lao, "Lag-net: Multi-granularity network for person re-identification via local attention system," *IEEE TMM*, 2021.

- [87] Lin Wu, Yang Wang, Junbin Gao, and Xue Li, “Where-and-when to look: Deep siamese attention networks for video-based person re-identification,” *IEEE TMM*, vol. 21, no. 6, pp. 1412–1424, 2018.
- [88] Yapkan Choi, Yeshwanth Napoleon, and Jan C van Gemert, “The arm-swing is discriminative in video gait recognition for athlete re-identification,” *arXiv preprint arXiv:2106.11280*, 2021.