

REFERENCE-AIDED PART-ALIGNED FEATURE DISENTANGLING FOR VIDEO PERSON RE-IDENTIFICATION

Guoqing Zhang, Yuhao Chen, Yang Dai, Yuhui Zheng, Yi Wu

ABSTRACT

Recently, video-based person re-identification (re-ID) has drawn increasing attention in compute vision community because of its practical application prospects. Due to the inaccurate person detections and pose changes, pedestrian misalignment significantly increases the difficulty of feature extraction and matching. To address this problem, in this paper, we propose a **Reference-Aided Part-Aligned (RAPA)** framework to disentangle robust features of different parts. Firstly, in order to obtain better references between different videos, a pose-based reference feature learning module is introduced. Secondly, an effective relation-based part feature disentangling module is explored to align frames within each video. By means of using both modules, the informative parts of pedestrian in videos are well aligned and more discriminative feature representation is generated. Comprehensive experiments on three widely-used benchmarks, i.e. iLIDS-VID, PRID-2011 and MARS datasets verify the effectiveness of the proposed framework. Our code will be made publicly available.

Index Terms— Person re-identification, Part alignment, Pose clues, Deep learning

1. INTRODUCTION

Person re-identification (re-ID) is an important retrieval task to match pedestrian images or videos captured from multiple non-overlapping cameras. Because of its wide application prospects in public safety and video surveillance, person re-ID has attracted increasing interest in recent years. Due to complicated and variable visual variations in practical scenarios such as pose, viewpoints, occlusion, illumination and background clutter, it remains a challenging task.

Currently, great progress has been made in image-based person re-ID [1, 2, 3], and video-based person re-ID has drawn increasing attention because of the impressive benefits of using multiple images, which can provide tremendous temporal information [4, 5]. In this paper, we focus on the person re-ID problem in the video setting.

Video-based re-ID task needs to aggregate features from multiple frames in video sequence. Some video-based person re-ID methods focus on global feature extraction and generating sequence-level features through traditional average pooling or spatial and temporal attention module [6, 7, 8, 9]. In

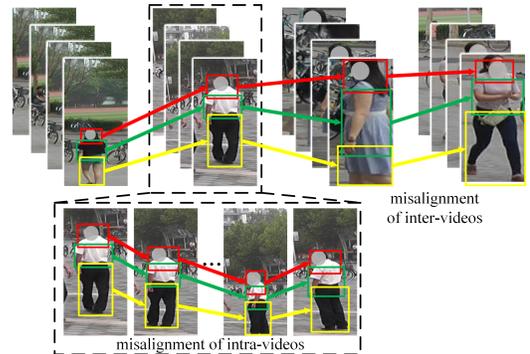


Fig. 1. The challenge of human region misalignment in video-based person re-ID.

some cases, extracting a global representation from the whole sequence may overlook local details. In order to mine local discriminative features, some methods divide global images into several body regions and learn global and local features simultaneously [1, 2, 5, 10]. However, because of the inaccurate person detections and pose changes, most part-based methods suffer from the problem of region misalignment of both intra- and inter-videos, as illustrated in Fig. 1. To align parts, some part-based methods take human parsing or semantic attributes into consideration and they will cost much more computation when a series of frames need to be preprocessed [3, 11]. To effectively deal with part misalignment and avoid excessive computation, we propose a **Reference-Aided Part-Aligned (RAPA)** framework for video-based person re-ID.

Our proposed RAPA framework is motivated by the application of reference-aided feature learning strategy [5] as well as the success of relation feature learning strategy [12, 13]. The architecture of our method mainly consists of a global feature extracting module, a reference feature learning module and a part feature disentangling module. Specifically, the global feature extracting module, which utilizes the global average pooling and temporal attention block, is applied to extract sequence-level features from the global point of view. The reference feature learning module is developed to find better reference frames and extract discriminative reference features. The part feature disentangling module is deployed to disentangle local features through aligning body parts of video sequences according to references.

We summarize the main contributions of our work into four aspects. First, we propose a novel **Reference-Aided Part-Aligned (RAPA)** framework for video based person re-ID, which aims to disentangle the discriminative features of different parts. Second, we develop a pose-based **Reference Feature Learning (RFL)** module to provide the uniform standard for alignment. Several discriminative reference features are extracted from reference frames to ensure the accurate alignment between different videos. Third, we design a relation-based **Part Feature Disentangling (PFD)** module, which aligns the body parts of intra-video. In this module, a relation-based attention block is adopted to search for corresponding body parts across frames. Finally, we evaluate the performance of the proposed RAPA framework on three mainstream benchmarks: MARS, iLIDS-VID and PRID-2011. Comprehensive experiments display show the superiority of our method to the state-of-the-arts.

2. RELATED WORK

Video-based Person Re-identification. Compared with image-based person re-ID, due to the more practical application prospects and much richer spatial-temporal information, video-based person re-ID has achieved more and more attention. Recurrent neural network (RNN) is a common network widely applied to analyze video sequence data. [6] and [8] introduced the models combining CNN and RNN to extract frame-level features and aggregate them by temporal pooling. However, these models treat all frames equally so that poor-quality frames and extraneous spatial regions influence the discriminability of features. To make the network focus on more informative frames and regions, attention mechanism is applied in video-based person re-identification. [9] proposed temporal attention models to calculate weight scores for time steps. Furthermore, [4] and [14] adopted attention mechanism in both spatial and temporal dimension. Compared with these existing video-based attention methods, our attention module in the proposed RAPA framework needs no additional parameters and mines attention scores according to the relations between reference features and frame-level feature maps.

Part-based Person Re-identification. In order to mine local informative details, some recent works divided pedestrian images into several parts and focused on local discriminative feature learning. One of the most intuitive and concise partition strategies is hard segmentation [1, 2]. [1] introduced a part-based convolutional baseline which divided pedestrian images vertically and equidistantly. [2] adopted a multiple granularities partition strategy to capture features of different scales. Compared with the hard segmentation, segmentation based on human pose and semantic information can avoid the problem of misalignment of body parts [3, 11]. However, most of traditional pose-based and semantics-based methods are not applicable to video-based person re-ID due to the huge and complex computation when they preprocess all frames

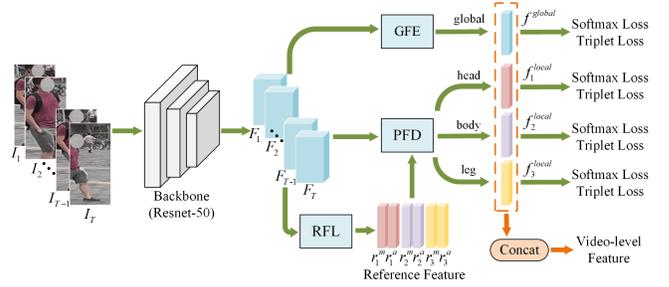


Fig. 2. The architecture of the proposed RAPA. This framework mainly contains three modules, including global feature extracting (GFE) module, reference feature learning (RFL) module and part feature disentangling (PFD) module. The GFE module is applied for the global spatiotemporal feature extraction, and RFL and PFD modules are jointly deployed to extract several local features.

of videos. Instead, our method only segments the reference frame according to pose information and utilizes references to align human parts.

Reference-aided Feature Learning Strategy. Owing to the continuity of video sequence, no significant difference exists in appearance or structure between consecutive frames. Accordingly, a single frame can become a reference and provide guidance for the analysis of the whole video sequence. Some recent works [5, 15] have applied reference-aided feature learning strategy in person re-ID. In [5], the first frame of each video is processed to guide subsequent frames, while in [15], the average of a video sequence is regarded as the reference frame. However, some poor-quality frames cannot provide enough reference value for video sequences. The quality evaluation block in our proposed RFL module solves this problem, which can select high quality frames automatically to generate better references.

3. PROPOSED METHOD

3.1. Framework Overview

The overview of our proposed framework is shown in Fig. 2. We select T frames from a video sequence as $V = \{I_t\}_{t=1}^T$, and the feature map of each frame $F_t \in \mathbb{R}^{C \times H \times W}$ is extracted through the backbone (e.g., ResNet-50), where C , H and W represent the channel, height and width of the feature maps, respectively. In the global branch, the feature maps are fed into Global Feature Extracting (GFE) module to extract features from the global point of view. Following the methods proposed in [5, 7], the feature maps are aggregated into image-level representations by the global average pooling, and then they are fused to video-level representations through temporal attention mechanism. After that, we use a 1×1 convolutional layer to reduce the dimension of features and get final global features denoted as $f^{global} \in \mathbb{R}^{\frac{C}{s}}$, where

s controls the dimension reduction ratio.

The local branch mainly contains Reference Feature Learning (RFL) module and Part Feature Disentangling (PFD) module. The former is used to find better local reference features, which can provide guidance for the alignment of video sequences, and the latter is used to align part features according to the references. Three local features (head, body and leg parts) extracted from both mentioned modules are denoted as $f_p^{local} \in \mathbb{R}^{\frac{C}{s}}$ ($p \in [1, 3]$). These two modules will be explained in detail in following subsections. The final video-based representation $f \in \mathbb{R}^{4 \times \frac{C}{s}}$ can be obtained by concatenating the global and local features:

$$f = [f^{global}, f_1^{local}, f_2^{local}, f_3^{local}] \quad (1)$$

3.2. Pose-based Reference Feature Learning Module

In [5], the first frame of the input video sequence is taken as a reference frame, which may be not in good condition due to the inaccurate person detections and occlusions. The quality of reference features determines the effect of alignment between videos. Therefore, we develop a posed-based reference feature learning (RFL) module to generate high-quality reference features and align part features between different videos, as illustrated in Fig. 3.

Firstly, in order to estimate the quality of images and find a better reference frame, we design a quality evaluation block which is motivated by temporal attention. Given the feature maps $F_t \in \mathbb{R}^{C \times H \times W}$ ($t \in [1, T]$), we get the image-level feature vectors $l_t \in \mathbb{R}^C$ with a global average pooling layer. The quality scores of frames are calculated by:

$$q_t = \text{Sigmoid}(\text{BN}(\text{Conv}(l_t))) \quad (2)$$

where the convolutional layer (Conv) reduces the vector dimension to 1 followed by a batch normalization layer (BN) and a sigmoid activation function (Sigmoid). The reference frame in each video sequence is defined as the frame which obtains the maximum quality score, denoted as I_k where $k = \arg \max(q_t)$.

Secondly, we apply the human pose estimation model (e.g., HRNet [16]) on I_k to predict human key-points. According to the distribution of key-points, the human body is divided into three parts including head, body and leg. The ROI pooling is commonly applied to capture the regions of interest from the whole feature map, which is widely used in object detection. In RFL module, both max and average ROI pooling are used to extract local reference features, because the former can focus on the image texture and the latter can ensure the integrity of information. We denote the local reference features as $r_1^m, r_1^a, r_2^m, r_2^a, r_3^m$ and $r_3^a \in \mathbb{R}^C$ (m means max pooling, a means average pooling, and 1-3 means head, body and leg), respectively.

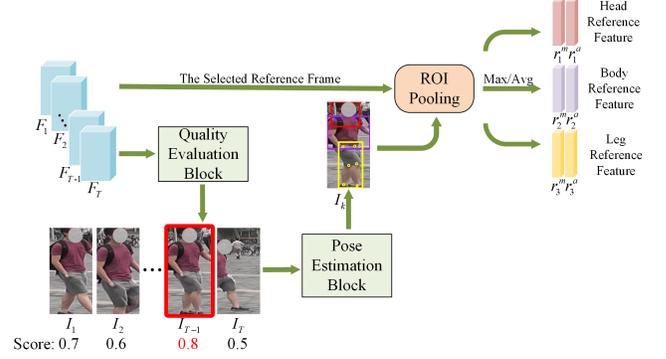


Fig. 3. The details of the pose-based Reference Feature Learning (RFL) module.

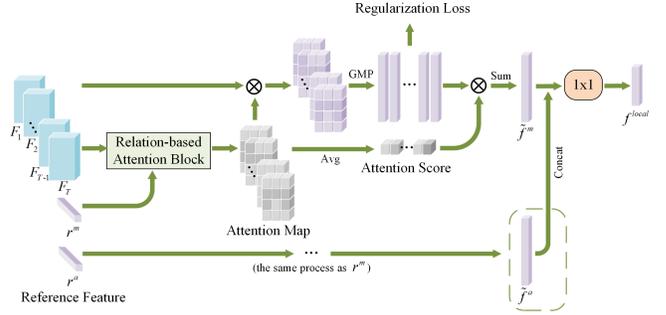


Fig. 4. The details of the relation-based Part Feature Disentangling (PFD) module.

3.3. Relation-based Part Feature Disentangling Module

The reference features provide the guidance for the alignment of intra-video sequences. To precisely disentangle local features, we introduce a relation-based part feature disentangling (PFD) module as shown in Fig. 4. In this paper, the vector along the channel dimension in feature maps is denoted as a column vector. We can measure the relevance between each reference feature vector and each column vector to obtain local attention maps in the relation-based attention block. Given the reference feature vectors r_p^m, r_p^a ($p \in [1, 3]$) and column vectors $v_t^{h,w} \in \mathbb{R}^C$ in feature maps F_t ($t \in [1, T]$, $h \in [1, H]$, $w \in [1, W]$), each relation element in relation map $D_{p,t}^m \in \mathbb{R}^{C \times H \times W}$ is calculated by:

$$d_{p,t,h,w}^m = \left(v_t^{h,w} - r_p^m \right)^2 \quad (3)$$

After that, a batch normalization layer and a sigmoid activation function are applied to normalize each relation element to the range of $(0, 1)$ and obtain the attention map $A_{p,t}^m \in \mathbb{R}^{C \times H \times W}$ by:

$$A_{p,t}^m = E\text{-Sigmoid}(\text{BN}(D_{p,t}^m)) \quad (4)$$

where $E \in \mathbb{R}^{C \times H \times W}$ is a matrix in which all elements are 1. Through the spatial attention mechanism and global max pooling (GMP), the elements in feature maps with high relevance can be found and aggregated to image-level local features $f_{p,t}^m \in \mathbb{R}^C$ as formulated:

$$f_{p,t}^m = \text{GMP}(F_t * A_{p,t}^m) \quad (5)$$

where $*$ is Hadamard product. Besides, in order to promote this module to focus on the more informative frames, the temporal channel attention mechanism is applied to weight the image-level local features. Based on the attention map, the attention score $S_{p,t}^m \in \mathbb{R}^C$ is computed as:

$$S_{p,t}^m = \frac{1}{H \times W} A_{p,t}^m \quad (6)$$

Then we can get the video-level aligned local feature $\tilde{f}_p^m \in \mathbb{R}^C$ through weighted sum:

$$\tilde{f}_p^m = \sum_{t=1}^T f_{p,t}^m * S_{p,t}^m \quad (7)$$

where $*$ is Hadamard product. The calculation of \tilde{f}_p^a is similar to \tilde{f}_p^m , and we omit the description of this part for convenience. Finally, we concatenate these two aligned features as $[\tilde{f}_p^m, \tilde{f}_p^a]$ and get the final part feature $f_p^{local} \in \mathbb{R}^{\frac{c}{s}}$ by performing a 1×1 convolutional layer on it to reduce its dimension, where s controls the dimension reduction ratio.

Due to the similarity of continuous frames, we design the inter-frame regularization term to promote the relation-based attention block to maintain the similarity between attention maps of different frames and avoid focusing on only one frame. Specifically, the regularization term of each video sequence is:

$$Reg = \sum_{i=1}^T \sum_{j=1, j \neq i}^T \sum_{p=1}^3 \left(\|f_{p,i}^m - f_{p,j}^m\|_2^2 + \|f_{p,i}^a - f_{p,j}^a\|_2^2 \right) \quad (8)$$

3.4. Loss Function

In our framework, we adopt both batch hard triplet loss and softmax cross entropy loss on each mini branch, as shown in Fig. 2. We assume that each mini-batch consists of P identities and K tracklets of each identity. The triplet loss for each branch is calculated by:

$$L_{tri} = \sum_{i=1}^P \sum_{a=1}^K \left[m + \overbrace{\max_{p=1 \dots K} \|f_a^{(i)} - f_p^{(i)}\|_2}^{\text{hardest positive}} - \underbrace{\min_{\substack{n=1 \dots K \\ j=1 \dots P \\ j \neq i}} \|f_a^{(i)} - f_n^{(j)}\|_2}_{\text{hardest negative}} \right]_+ \quad (9)$$

where $f_a^{(i)}$, $f_p^{(i)}$ and $f_n^{(j)}$ are the features extracted from the anchor, positive and negative samples respectively, and m is the margin hyperparameter to control the differences between intra and inter distances. The softmax cross entropy loss for each branch is formulated as:

$$L_{softmax} = -\frac{1}{P \times K} \sum_{i=1}^P \sum_{a=1}^K y_{i,a} \log q_{i,a} \quad (10)$$

where $y_{i,a}$ and $q_{i,a}$ are the ground truth identity and prediction of tracklet sample $\{i, a\}$. Therefore, the total loss for each branch is:

$$L_c^{branch} = L_{tri} + L_{softmax} \quad (11)$$

where $c \in [1, 4]$ indicates the branch number and L_1^{branch} and $\{L_c^{branch}\}_{c=2}^4$ indicate the loss from global branch and three local branches, respectively. Besides, the inter-frame regularization loss can be calculated by:

$$L_{reg} = \frac{1}{P \times K} \sum_{i=1}^P \sum_{a=1}^K Reg \quad (12)$$

The addition of branch losses and regularization loss constitutes the final loss for optimization:

$$L_{total} = \sum_{c=1}^4 L_c^{branch} + \lambda L_{reg} \quad (13)$$

where λ is a hyper-parameter to control the proportion of regularization loss.

4. EXPERIMENTS

4.1. Datasets and Evaluation Protocol

Datasets. Three standard video-based person re-ID datasets: iLIDS-VID [17], PRID-2011 [18] and MARS [19], are applied to evaluate our proposed framework. **iLIDS-VID** dataset, which is challenging due to occlusion and blur, consists of 300 identities and each identity includes 2 video sequences taken from a pair of non-overlapping cameras. By comparison, **PRID-2011** dataset is less challenging due to its simple environments with rare occlusion. It contains 385 and 749 identities from 2 different cameras, but only the first 200 identities appear in both cameras. **MARS** dataset is one of the largest video-based person re-ID benchmarks, which consists of 1261 identities and 20,715 video sequences captured by 6 cameras.

Evaluation Protocol. The CMC curve and mAP are applied to evaluate the performance of our framework. For iLIDS-VID and PRID-2011 datasets, following the common practices in previous work [17], we randomly split the half-half for training and testing. The average result of 10 repeated experiments is reported. For MARS dataset, 8298 sequences of 625 identities are used for training and other sequences are used for testing.

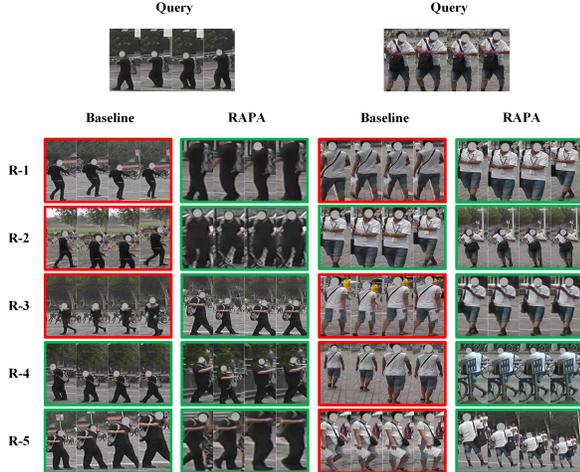


Fig. 5. Visualization of person re-ID results using the baseline model and our proposed RAPA framework. The green and red bounding boxes indicate correct and incorrect matches, respectively.

Table 1. Ablation study on the components of our proposed method on MARS dataset.

Variant	Rank-1	Rank-5	Rank-20	mAP
(a) Baseline	82.4	93.8	97.1	74.1
(b) Baseline+RA	84.7	95.1	98.1	76.3
(c) Baseline+RA+Reg	86.7	96.2	98.1	81.0
(d) Baseline+RA+Reg+TA	87.7	96.1	98.2	82.2
(e) Baseline+RA+Reg+TA+PE	88.0	96.5	98.2	82.7
(f) Baseline+RA+Reg+TA+PE+QE	88.7	96.1	98.1	82.8

4.2. Implementation Details

In the training phase, we randomly select $T = 4$ frames of each video as the input sequence. Input images are resized to 128×256 . Random erasing and random cropping are applied for data augmentation. A mini-batch consists of 8 identities with 4 tracklets. The ResNet-50 pretrained on ImageNet is utilized as our backbone. We choose Adam to optimize our model with weight decay 5×10^{-4} . The initial learning rate is 3.5×10^{-4} and decreased by 0.1 every 100 epochs. The epoch number is 400 in total. For iLIDS-VID and PRID-2011 datasets, the hyper-parameter λ in final loss function is set to 5×10^{-5} , and for MARS dataset, λ is set to 3×10^{-4} . In the testing phase, the video sequence is segmented into several clips of length $T = 4$. The average of all clip-level features of the same sequence is regarded as the video-level feature. Euclidean distance is applied to measure the similarity between query sequences and gallery sequences.

4.3. Ablation Study

To evaluate the effectiveness of components in our proposed RAPA framework, we conduct a series of ablation experiments and show the comparative results in Table 1. We se-

Table 2. Ablation study on the branches of feature representation on MARS dataset.

Variant	Rank-1	Rank-5	Rank-20	mAP
(a) Global (Baseline)	82.4	93.8	97.1	74.1
(b) Local	87.4	96.6	98.3	81.7
(c) Global+Local	88.7	96.1	98.1	82.8

lect the global branch without 1×1 convolutional layer as our baseline, which follows the method proposed in [7]. The PFD module is divided into **RA**, **Reg** and **TA**, which correspond to the relation-based attention block, the inter-frame regularization loss and the temporal channel attention score, respectively. The RFL module includes **PE** and **QE**, which indicate the pose estimation block and the quality evaluation block, respectively. Compared with the baseline, our framework improves Rank-1 and mAP from 82.4% and 74.1% to **88.7%** and **82.8%** on MARS dataset. Some comparative results are visualized in Fig. 5. As can be observed intuitively, misalignments bring difficulties for the baseline model to distinguish some pedestrians with similar appearances, while our proposed RAPA framework can achieve more robust results in these cases.

Effectiveness of PFD Module. The comparative results of variants (a) – (f) show the effectiveness of our proposed PFD module. Variant (a) doesn’t perform well because it only takes the global feature into consideration but ignores the more discriminative local details. Variant (b) utilizes the relation-based attention block which can focus on the corresponding local areas. Without the RFL module, variant (b) applies the hard segmentation on the first frame to generate the reference features by default. It can be observed that **RA** improves Rank-1 and mAP accuracy by 2.3% and 2.2%, respectively. The application of **Reg** in variant (c) preserves the similarity between continuous frames and further improves the accuracy of our framework. To encourage the framework to focus on the frames of interest, **TA** is adopted in variant (d) which forms a complete PFD module. In summary, compared with the baseline, our framework with PFD module achieves 5.3% and 8.1% improvements in Rank-1 and mAP, respectively.

Effectiveness of RFL Module. The comparable results of variants (e) and (f) prove the effectiveness of our proposed RFL module. Variant (e) utilizes **PE** block to segment frames into several parts according to the pose information. Compared with variant (d), variant (e) solves the misalignment between different video sequences. Variant (f) further deploys **QE** block to find high-quality reference frames which improves the robustness in some complex cases such as occlusion. Finally, variant (f) can outperform variant (d) by 1% and 0.6% in Rank-1 and mAP on MARS dataset.

Effectiveness of Different Branches. As shown in Fig. 2, the feature representation in our framework consists of global branch and 3 local branches. To verify the effective-

Table 3. Comparisons of our proposed method to the state-of-the-art methods on MARS, iLIDS-VID and PRID-2011 datasets. The 1st, 2nd and 3rd best results are emphasized with red, blue and green color, respectively.

Method	Publication	MARS		iLIDS-VID	PRID-2011
		Rank-1	mAP	Rank-1	Rank-1
RNN [6]	CVPR'16	-	-	58.0	70.0
CNN+XQDA [19]	ECCV'16	68.3	49.3	53.0	77.3
CRF [8]	CVPR'17	71.0	-	61.0	77.0
SeeForest [9]	CVPR'17	70.6	50.7	55.2	79.4
RQEN [20]	AAAI'18	77.8	71.1	76.1	92.4
STAN [14]	CVPR'18	82.3	65.8	80.2	91.2
STA [4]	AAAI'19	86.3	80.1	-	-
RRU [21]	AAAI'19	84.4	72.7	84.3	92.7
A3D [22]	TIP'20	86.3	80.4	87.9	95.1
AMEM [23]	AAAI'20	86.7	79.3	87.2	93.3
FGRA [5]	AAAI'20	87.3	81.2	88.0	95.5
Two-stream M3D [24]	TIP'20	88.6	79.5	86.7	96.6
Ours	ICME'21	88.7	82.8	89.6	95.2

ness of these branches, several ablation experiments are conducted and the comparable results are shown in Table 2. Variant (a) only contains the global branch which is our baseline. Variant (b) disentangles several local features and achieves a significant improvement. Compared with variant (b), variant (c) achieves the best results owing to the information compensation from global branch.

4.4. Comparison with the State-of-the-arts

Table 3 reports the results of our proposed RAPA framework and some state-of-the-art methods on MARS, iLIDS-VID and PRID-2011 datasets. As we can observe from Table 3, our model obtains the Rank-1 accuracy of **88.7%**, **89.6%** and **95.2%** on MARS, iLIDS-VID and PRID-2011, which outperforms all the state-of-the-art methods in most evaluation indicators. The main reason is that our framework accomplishes the feature alignment of both intra- and inter-video sequences. However, our RAPA framework has the slightly lower accuracy than FGRA and Two-stream M3D on PRID-2011 dataset. One possible explanation is that the pedestrian images in PRID-2011 are perfectly neat and the condition of misalignment is rare, so that the alignment strategy cannot exert its advantage.

5. CONCLUSION

In this paper, a reference-aided part-aligned (RAPA) framework is proposed for video-based person re-identification. In order to solve the pedestrian misalignment problem, the pose-based reference feature learning (RFL) module and the relation-based part feature disentangling (PFD) module are explored. The former is designed to extract discriminative reference features and align the local features between different video sequences. The latter is applied to align parts of intra-video sequences according to the references. Moreover, in PFD module, a relation-based attention block is adopted to

search for corresponding body parts. The outstanding experimental results prove the superiority of the proposed RAPA framework.

6. REFERENCES

- [1] Sun et al., "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," in *ECCV*, 2018.
- [2] Wang et al., "Learning discriminative features with multiple granularities for person re-identification," in *ACMMM*, 2018.
- [3] Jin et al., "Semantics-aligned representation learning for person re-identification.," in *AAAI*, 2020.
- [4] Fu et al., "Sta: Spatial-temporal attention for large-scale video-based person re-identification," in *AAAI*, 2019.
- [5] Chen et al., "Frame-guided region-aligned representation for video person re-identification.," in *AAAI*, 2020.
- [6] McLaughlin et al., "Recurrent convolutional network for video-based person re-identification," in *CVPR*, 2016.
- [7] Gao et al., "Revisiting temporal modeling for video-based person reid.," *arXiv*, 2018.
- [8] Chen et al., "Deep spatial-temporal fusion network for video-based person re-identification," in *CVPR*, 2017.
- [9] Zhou et al., "See the forest for the trees: Joint spatial and temporal recurrent neural networks for video-based person reidentification," in *CVPR*, 2017.
- [10] Sun et al., "Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification," in *CVPR*, 2019.
- [11] Zhao et al., "Spindle net: Person re-identification with human body region guided feature decomposition and fusion," in *CVPR*, 2017.
- [12] Zhang et al., "Relation-aware global attention for person re-identification," in *CVPR*, 2020.
- [13] Park et al., "Relation network for person re-identification," *arXiv*, 2019.
- [14] Li et al., "Diversity regularized spatiotemporal attention for video-based person re-identification," in *CVPR*, 2018.
- [15] Zhang et al., "Multi-granularity reference-aided attentive feature aggregation for video-based person re-identification," in *CVPR*, 2020.
- [16] Sun et al., "Deep high-resolution representation learning for human pose estimation," in *CVPR*, 2019.
- [17] Wang et al., "Person re-identification by video ranking," in *ECCV*, 2014.
- [18] Hirzer et al., "Person re-identification by descriptive and discriminative classification," in *SCIA*, 2011.
- [19] Zheng et al., "Mars: A video benchmark for large-scale person re-identification," in *ECCV*, 2016.
- [20] Song et al., "Region-based quality estimation network for large-scale person re-identification," in *AAAI*, 2018.
- [21] Liu et al., "Spatial and temporal mutual promotion for video-based person re-identification," in *AAAI*, 2019.
- [22] Chen et al., "Learning recurrent 3d attention for video-based person re-identification," *TIP*, 2020.
- [23] Li et al., "Appearance and motion enhancement for video-based person re-identification," in *AAAI*, 2020.
- [24] Li et al., "Multi-scale temporal cues learning for video person re-identification," *TIP*, 2020