

INTRA-CLIP AGGREGATION FOR VIDEO PERSON RE-IDENTIFICATION

Takashi Isobe[†] Jian Han[†] Fang Zhu[‡] Yali Li[†] Shengjin Wang^{†✉}

[†]Beijing National Research Center for Information Science and Technology,
Department of Electronic Engineering, Tsinghua University, Beijing 100084, China

[‡]MetroTech Center, Brooklyn, New York University, 11201, USA

ABSTRACT

Video-based person re-identification has drawn massive attention in recent years due to its extensive applications in video surveillance. While deep learning-based methods have led to significant progress, these methods are limited by ineffectively using complementary information, which is blamed on necessary data augmentation in the training process. Data augmentation has been widely used to mitigate the over-fitting trap and improve the ability of network representation. However, the previous methods adopt image-based data augmentation scheme to individually process the input frames, which corrupts the complementary information between consecutive frames and causes performance degradation. Extensive experiments on three benchmark datasets demonstrate that our framework outperforms the most recent state-of-the-art methods. We also perform cross-dataset validation to prove the generality of our method.

Index Terms— Video Person Re-identification, Deep Learning, Data Augmentation.

1. INTRODUCTION

Person re-identification (re-ID) aims to recognize the same identity in different images or videos captured by different cameras distributed at separated physical locations. In contrast with image person re-ID [1, 2, 3, 4], video person re-ID is more robust to noise. Both spatial information across positions and temporal information across frames can be used to represent clip-level features. The previous works [5, 6] exploit motion estimation either implicitly (*e.g.* gait) or explicitly (*e.g.* optical flow) to represent a video sequence. But those works are not optimal for video person re-ID. Inaccurate motion estimation, especially when there is occlusion or parallax, deteriorating the final performance. Besides, those methods often suffer a heavy computational load. To represent a clip-level descriptor while maintaining a low computational cost, temporal pooling has been widely used in recent works [9, 10]. They perform generic or weighted average pooling in the end of the network to aggregate intra-clip features across time. However, temporal pooling is a linear operation which is limited to capture the specific features of

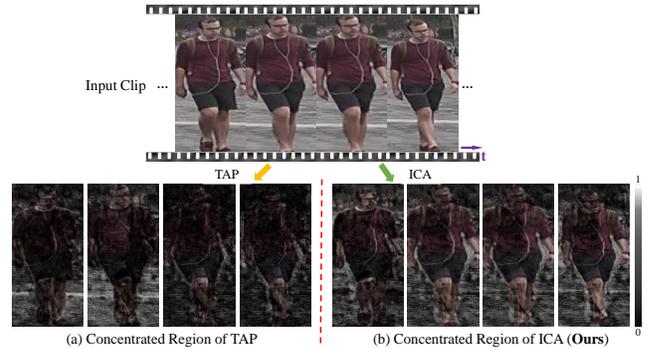


Fig. 1. Visualization of the concentrated regions for the intra-clip frames. For fair comparison, we use same feature extraction backbone, *i.e.* ResNet50 [7], but with different temporal aggregation strategy. (a) and (b) adopt Temporal Average Pooling (TAP) and the proposed hierarchical aggregation strategy, respectively. In (a), the activated maps have a scattered distribution with less meaningfulness. In (b), the activated maps are more concentrated and meaningful around body parts. The activated maps are obtained by Grad-cam [8].

a video sequence. In this paper, we propose a novel Intra-Clip Aggregation (ICA) module in order to effectively integrate the clip-level features. Specifically, ICA is a cascade structure which consists of a learnable block followed by a temporal average pooling layer. The critical innovation of the ICA is hierarchically aggregating the intra-clip features with both linear and non-linear operations. The linear operation is used to generate global features, and then we apply a non-linear block to describe the most important semantic concepts of clip-level features.

Data augmentation, an explicit form of regularization, has been widely used in the training process of deep neural network. The general data augmentation approaches such as random cropping, flipping as well as erasing [11, 12] work well on image person re-ID task by randomly transferring or noising the original images. The previous works [9, 13, 6] treat video person re-ID as a generalized image person re-ID task. They apply data augmentation operation asynchronously on the input frames. In such process, each frame

is transformed with a random probability, which introduces excessive noise and consequently corrupts the temporal cues. For example, randomly flipping each frame will result in misalignment. In addition, randomly erasing a region of pixels of each frame will cause too much spatial-temporal information loss. In this case, the model may be confused to fully utilize the interactive information to represent informative clip-level features and perform poorly in the presence of real-world noise, *e.g.*, occlusion, lighting and motion blur. In this paper, we propose a video-based data augmentation approach for video person re-ID, which is a temporal extension of commonly used image-level data augmentation techniques. The proposed video-based data augmentation is easy to implement and meanwhile yields consistent improvement over three challenging video person re-ID benchmarks.

To sum up, our contributions are listed as following:

- We revisit data augmentation for video person re-ID task, and propose a novel video-based data augmentation strategy to strengthen the representation ability and the generality of the learned model. It can be adopted on various existed image-based data augmentation approaches.
- We propose a novel cascade temporal integration pipeline which effectively integrates the intra-clip features in a hierarchical manner.
- Our model outperforms the state-of-the-art methods on ILIDS-VID [14] and MARS [15] benchmarks by a large margin. The impressive result on cross-data validation shows the generality of the proposed method.

2. METHODOLOGY

In training process, the intra-clip frames within a mini-batch either are augmented synchronously or remain unchanged, and then a base network is applied to extract the frame-level features. Finally, ICA module deeply integrates intra-clip features in a hierarchical manner to represent the clip-level pedestrian descriptors. Fig. 2 shows the proposed ICA module. The main contributions of the overall framework are two parts: 1) video-based data augmentation; 2) hierarchically temporal integration module. More details about synchronous data augmentation and ICA module are presented in Sec.2.1 and Sec.2.2, respectively.

2.1. Synchronous Data Augmentation

In this subsection, we improve the data augmentation strategy for video person re-ID task. Commonly used image-level data augmentation [11, 12] approaches perform well on image person re-ID by suppressing the underlying noise such as camera intrinsic noise and background noise. However, as for a video sequence, asynchronous data augmentation

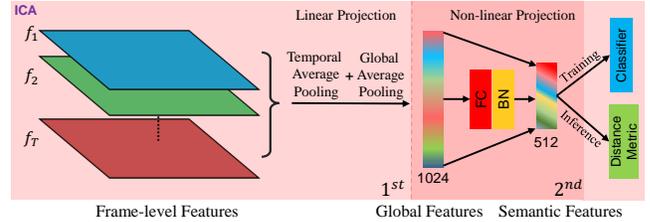


Fig. 2. The proposed ICA module which integrates the intra-clip information with a hierarchical pattern.

may introduce unnecessary noise corrupting the temporal cues of intra-clip frames, which leads to the result that the model poorly resists the noise from the real world and may be confused about how to utilize the intra-clip complementary information. To address aforementioned drawbacks, we propose a novel video-based data augmentation strategy, termed as Synchronous Data Augmentation (SDA). Our method can effectively preserve the complementary information among consecutive frames and change the underlying noise among frames synchronously, which helps the network to learn a discriminative distance metric and better utilize the interactive information among frames. In training, the intra-clip frames within a mini-batch randomly undergo either of the two operations: 1) remaining unchanged; 2) being synchronously transformed with commonly used data augmentation techniques such as random flipping [11] and random erasing [12]. We formulate the operation of asynchronous transformation and the proposed asynchronous transformation as following. For simplicity, we formulate one data augmentation process as example.

asynchronous transformation:

$$T_{at}\{f_1, f_2, \dots, f_T\} = \{\psi_k(f_k)\}_{k=1}^T \quad (1)$$

synchronous transformation:

$$T_{st}\{f_1, f_2, \dots, f_T\} = \{\psi(f_k)\}_{k=1}^T \quad (2)$$

where $\Psi(\cdot)$ denotes the operator of data augmentation. For asynchronous data augmentation $T_{at}\{\cdot\}$, the operation $\{\Psi_k(\cdot)\}_{k=1}^T$ is randomly changed over the input T frames. For synchronous data augmentation $T_{st}\{\cdot\}$, all frames are applied with the same operation $\Psi(\cdot)$. $\{f_k\}_{k=1}^T$ denotes T frames of a tracklet.

In this paper, we incorporate three types of augmentation approaches, *i.e.*, random cropping, flipping and erasing. The transformation probability of each operation is fixed along temporal axis, *i.e.*, cropping size, rotating angle and erasing region.

2.2. Intra-clip Aggregation Module

The key idea of ICA module is to capture the important semantic concepts of clip-level features. Based on such motivation, ICA is designed as a cascade structure, which can fully

integrate clip-level information in a hierarchical manner. ICA takes the frame-level features as input, and performs average pooling to generate clip-level global features in preliminary fusion, which can be expressed as:

$$Z_{1,c,1,1} = \frac{1}{WHT} \sum_{w=1}^W \sum_{h=1}^H \sum_{t=1}^T X_{t,c,w,h} \quad (3)$$

Where $X_{t,c,w,h}$ is temporally concatenated frame-level features. $Z_{1,c,1,1}$ is clip-level global features obtained by linear projection. Subsequently, the above $Z_{1,c,1,1}$ is further integrated with a high dimensional feature projection block, which can be formulated as:

$$Y_{\tilde{c}} = G_{c \rightarrow \tilde{c}}(Z_{1,c,1,1}) \quad (4)$$

Where $Y_{\tilde{c}}$ stands for the clip-level semantic embedding. $G_{c \rightarrow \tilde{c}}(\cdot)$ denotes the non-linear projection. To reduce the computational cost, we design $G_{c \rightarrow \tilde{c}}(\cdot)$ as a bottleneck structure, which is composed of fully-connected (FC) layer and batch normalization (BN) [16]. More details about the parameter setting and structure of ICA can be found in Fig. 2.

In comparison with the most existing methods [9, 10], which typically adopt a temporal pooling layer in the end of the network to represent the clip-level features, our method is able to generate more discriminative clip-level features by leveraging the interactional information among consecutive frames (see Fig. 1). The impressive experimental results about synchronous data augmentation and ICA module are presented in Sec.3.

3. EXPERIMENTS

3.1. Datasets and Evaluation Protocols

We conduct extensive experiments on three challenging video-based person re-ID datasets, including the ILIDS-VID [14], PRID 2011 [17] and MARS [15]. **The ILIDS-VID dataset** contains 300 identities forming a total number of 600 video sequences. Pedestrians are observed in two non-overlapping cameras views. The length of each tracklet varies from 23 to 192, with an average number of 73 frames. The train and test set are splitted evenly of 150 identities. **The PRID 2011 dataset** is another standard benchmark for video-based person re-ID. Following [18, 19], we select 178 out of 200 identities with more than 21 frames forming a total number of 356 video sequences. **The MARS dataset** is the one of the largest datasets, which consists of 1,261 different IDs and around 20,000 tracklets from 6 cameras. We evaluate the performance of the proposed method on ILIDS-VID and MARS, and implement cross-dataset evaluation on PRID2011.

We evaluate our model by Mean Average Precision (mAP) score and Cumulative Matching Characteristic (CMC) curve.

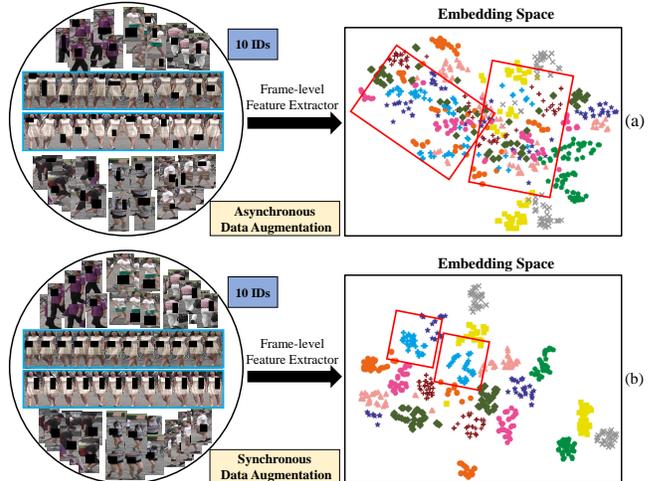


Fig. 3. Visualization of the embedding space. We sampled ten identities, and each ID contains two video clips with twelve frames. We use ResNet50, pre-trained on ImageNet, as frame-level feature extractor, and then visualize these features by t-SNE [20].

Datasets Rank@k	Model	ILIDS-VID			MARS			
		1	5	20	1	5	20	mAP
Baseline	Model 1	78.2	87.8	92.1	77.8	86.4	90.7	74.5
Baseline + SDA	Model 2	80.3	89.5	92.8	79.1	87.4	91.4	75.6
ICA	Model 3	86.1	96.9	98.7	85.1	95.2	97.4	80.7
ICA + SDA	Model 4	88.7	98.7	100.0	87.5	96.6	98.2	81.6

Table 1. Ablation on the effectiveness of the proposed components.

3.2. Implementation Details

Our model is supervised by triplet loss [21] and cross-entropy loss [22]. Adam [23] optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ is adopted to optimize the proposed framework, where weight decay is set to 5×10^{-4} . The learning rate is initially set to 4×10^{-4} and later down-scaled by a factor of $\frac{1}{200}$ every 200 epochs until 500 epochs. We randomly sample 4 identities with 8 consecutive clips forming the mini-batch of 32. According to [9, 24], we set the length of each clip to 4. The input frames are uniformly resized to $256 \times 128 \times 3$ and linearly scaled to $[-1, 1]$. In training, the temporally random cropping, temporally random flipping and temporally random erasing are used to augment each video clip. We adopt the same settings for training all datasets. We utilize the ResNet50 [7] pre-trained on ImageNet as the frame-level feature extraction network. Note that our method can be easily generalized to other backbones. All experiments are conducted on a server with Python 3.6.4 and Pytorch 1.1 platform.

Datasets Rank@k	Backbone	ILIDS-VID			MARS			
		1	5	20	1	5	20	mAP
Baseline + SDA	Alexnet [11]	53.1	69.4	75.1	51.3	66.8	73.8	49.7
	InceptionV3 [25]	69.5	83.6	89.7	67.2	81.6	87.3	62.5
	ResNet50 [7]	80.3	89.5	92.8	79.1	87.4	91.4	75.6
ICA + SDA	Alexnet [11]	62.1	79.0	83.1	59.5	76.1	80.3	54.7
	InceptionV3 [25]	78.4	94.1	97.4	76.3	91.4	95.0	72.0
	ResNet50 [7]	88.7	98.7	100.0	87.5	96.6	98.2	81.6

Table 2. Ablation on the effectiveness of different feature extractors.

3.3. Ablation Study

The baseline (Model 1) corresponds to resnet-50 with temporal average pooling, and training with asynchronously random cropping, flipping and erasing. The effectiveness of each component is reported in Table 1. From top to bottom, we evaluate each component successively. We can observe that, with the proposed video-based data augmentation, Model 2 surpasses the Model 1 on ILIDS-VID and MARS datasets about 2.1% and 1.3% at rank-1, respectively. Attributed to ICA module, Model 3 improves the rank-1 accuracy by 7.9% and 7.3% than Model 1 on ILIDS-VID and MARS, respectively. By combining the proposed SDA and ICA module, Model 4 works better than other models by a large margin. We also implement our method on different backbones, and the result shows that “ICA+SDA” consistently outperforms baseline, demonstrating that our method performs well with different frame feature extractors.

To better understand the difference between asynchronous and synchronous data augmentation, we carefully visualize the corresponding embedding space, as shown in Fig. 3. We can see that the (a) illustrates a scattered distribution, where the intra-clip frames are separated and mixed with other clips. The (b) shows intra-clip frames are clustered and preserve more clip-level information.

3.4. Comparison with State-of-the-arts

We use the best model (Model 4) obtained by the proposed ICA and SDA to compare with previous state-of-the-art results. As shown in Table 3, the first two methods [6, 26] explicitly utilize the temporal information by using a on-line network to estimate the optical flow between the consecutive frames and combining them with the spatial information to represent the clip-level features. The third method [5] contributes to implicitly use the spatio-temporal information with a succession of 3D convolutions. The last seven methods [18, 13, 10, 9, 19, 24, 27] aggregate intra-clip features over temporal dimension to represent the clip-level features. **Note:** STMP [19] uses inceptionV3 as their backbone network. Different feature extractors have significantly impact on the final performance, as shown in Table 2. We carefully re-implement [19] by adopting ResNet50 as the backbone network. Table 3 reveals that our method outperforms other state-of-the-art methods by a large margin, more than

Datasets Rank@k	ILIDS-VID			MARS			
	1	5	20	1	5	20	mAP
QAN [6]	68.0	86.8	97.4	73.7	84.9	91.6	51.7
AMOC [26]	68.7	94.3	99.3	68.3	81.4	90.6	52.9
Liao <i>et al.</i> [5]	81.3	-	-	84.3	-	-	77.0
TRL [18]	57.7	81.7	94.1	79.3	91.1	96.0	66.8
STSRN [13]	70.0	89.3	98.7	76.7	93.8	98.1	-
DRSA [10]	80.2	-	-	82.3	-	-	65.8
Gao <i>et al.</i> [9]	-	-	-	83.3	93.8	97.4	76.7
STMP [19]	85.7	97.5	99.8	86.2	95.3	97.8	75.6
STA [24]	-	-	-	86.3	95.7	97.1	80.8
GLTR [27]	86.0	98.0	-	87.0	95.8	98.2	78.5
Ours	88.7	98.7	100.0	87.5	96.6	98.2	81.6

Table 3. Performance comparison with other state-of-the-art methods on ILIDS-VID and MARS datasets. “-”: no reported results.

Datasets Rank@k	PRID 2011		
	1	5	20
TRL [18]	29.5	59.4	82.2
STSRN [13]	30.0	58.0	85.0
STMP [19]	32.0	58.0	90.0
Ours	41.6	71.9	92.1

Table 4. Generality comparison with other state-of-the-arts methods.

0.5% better than the recent proposed STA [24] and GLTR [27] at rank-1 and rank-5 on MARS benchmark. Due to the overfitting trap, the previous methods may exhibit good performance on the training set, but suffers severe performance degeneration when work on a new dataset. To better understand the generalization performance of our method, we conducted cross-dataset experiment. We trained Model 4 on ILIDS-VID and evaluate it on PRID2011. Table 4 shows that our model achieves consistently superior performance over other methods, which demonstrates the generality of our method.

4. CONCLUSIONS

In this paper, we revisit data augmentation for video person re-ID task, and propose a video-based data augmentation scheme, termed as Synchronous Data Augmentation, for training the convolutional neural network. Benefited from the proposed data augmentation strategy, our model is better to utilize the interactive information among frames and has strong generality. In order to extract clip-level semantic features, we also propose a ICA module to integrate the intra-clip features in hierarchical manner. Thanks to the proposed data augmentation strategy and temporal integration pipeline, we achieve new state of the art on ILIDS-VID and MARS benchmarks without re-ranking. We also perform the cross-dataset validation and confirm the generality of our method.

5. REFERENCES

- [1] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 480–496.
- [2] Yifan Sun, Liang Zheng, Weijian Deng, and Shengjin Wang, “Svdnet for pedestrian retrieval,” in *ICCV*, 2017.
- [3] Wei Li, Xiatian Zhu, and Shaogang Gong, “Harmonious attention network for person re-identification,” in *CVPR*, 2018.
- [4] Takashi Isobe, Dong Li, Lu Tian, Weihua Chen, Yi Shan, and Shengjin Wang, “Towards discriminative representation learning for unsupervised person re-identification,” *arXiv preprint arXiv:2108.03439*, 2021.
- [5] Xingyu Liao, Lingxiao He, and Zhouwang Yang, “Video-based person re-identification via 3d convolutional networks and non-local attention,” *arXiv preprint arXiv:1807.05073*, 2018.
- [6] Yu Liu, Junjie Yan, and Wanli Ouyang, “Quality aware network for set to set recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5790–5799.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [8] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [9] Jiyang Gao and Ram Nevatia, “Revisiting temporal modeling for video-based person reid,” *arXiv preprint arXiv:1805.02104*, 2018.
- [10] Shuang Li, Slawomir Bak, Peter Carr, and Xiaogang Wang, “Diversity regularized spatiotemporal attention for video-based person re-identification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 369–378.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [12] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang, “Random erasing data augmentation,” *arXiv preprint arXiv:1708.04896*, 2017.
- [13] Xinxing Su, Yingtian Zou, Yu Cheng, Shuangjie Xu, Mo Yu, and Pan Zhou, “Spatial-temporal synergic residual learning for video person re-identification,” *arXiv preprint arXiv:1807.05799*, 2018.
- [14] Taiqing Wang, Shaogang Gong, Xiatian Zhu, and Shengjin Wang, “Person re-identification by video ranking,” in *European Conference on Computer Vision*. Springer, 2014, pp. 688–703.
- [15] 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 *European Conference on Computer Vision*. Springer, 2016, pp. 868–884.
- [16] Sergey Ioffe and Christian Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
- [17] Martin Hirzer, Csaba Beleznai, Peter M Roth, and Horst Bischof, “Person re-identification by descriptive and discriminative classification,” in *Scandinavian conference on Image analysis*. Springer, 2011, pp. 91–102.
- [18] Ju Dai, Pingping Zhang, Dong Wang, Huchuan Lu, and Hongyu Wang, “Video person re-identification by temporal residual learning,” *IEEE Transactions on Image Processing*, vol. 28, no. 3, pp. 1366–1377, 2019.
- [19] Yiheng Liu, Zhenxun Yuan, Wengang Zhou, and Houqiang Li, “Spatial and temporal mutual promotion for video-based person re-identification,” *arXiv preprint arXiv:1812.10305*, 2018.
- [20] Laurens van der Maaten and Geoffrey Hinton, “Visualizing data using t-sne,” *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [21] Alexander Hermans, Lucas Beyer, and Bastian Leibe, “In defense of the triplet loss for person re-identification,” *arXiv preprint arXiv:1703.07737*, 2017.
- [22] Christopher M Bishop, *Pattern recognition and machine learning*, springer, 2006.
- [23] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [24] Yang Fu, Xiaoyang Wang, Yunchao Wei, and Thomas Huang, “Sta: Spatial-temporal attention for large-scale video-based person re-identification,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 8287–8294.

- [25] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [26] Hao Liu, Zequn Jie, Karlekar Jayashree, Meibin Qi, Jianguo Jiang, Shuicheng Yan, and Jiashi Feng, “Video-based person re-identification with accumulative motion context,” *IEEE transactions on circuits and systems for video technology*, vol. 28, no. 10, pp. 2788–2802, 2018.
- [27] Jianing Li, Jingdong Wang, Qi Tian, Wen Gao, and Shiliang Zhang, “Global-local temporal representations for video person re-identification,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 3958–3967.