

GAITAKE: GAIT RECOGNITION BY TEMPORAL ATTENTION AND KEYPOINT-GUIDED EMBEDDING

Hung-Min Hsu¹, Yizhou Wang¹, Cheng-Yen Yang¹, Jenq-Neng Hwang¹, Hoang Le Uyen Thuc², Kwang-Ju Kim³

¹ Information Processing Lab, University of Washington, Seattle, USA

² University of Science and Technology, University of Danang, Danang, Vietnam

³ Electronics and Telecommunications Research Institute (ETRI), Daegu, South Korea

ABSTRACT

Gait recognition, which refers to the recognition or identification of a person based on their body shape and walking styles, derived from video data captured from a distance, is widely used in crime prevention, forensic identification, and social security. However, to the best of our knowledge, most of the existing methods use appearance, posture and temporal features without considering a learned temporal attention mechanism for global and local information fusion. In this paper, we propose a novel gait recognition framework, called Temporal Attention and Keypoint-guided Embedding (GaitTAKE), which effectively fuses temporal-attention-based global and local appearance feature and temporal aggregated human pose feature. Experimental results show that our proposed method achieves a new SOTA in gait recognition with rank-1 accuracy of 98.0% (normal), 97.5% (bag) and 92.2% (coat) on the CASIA-B gait dataset; 90.4% accuracy on the OU-MVLP gait dataset.

Index Terms— Gait Recognition, Temporal Attention, Human Pose Estimation

1. INTRODUCTION

Gait recognition, which uses video data captured from a distance to recognize or identify a person based on their body shape and walking styles, is widely used in crime prevention, forensic identification, and social security, etc. Person re-identification (ReID) is one of the most popular research in the computer vision community. However, merely using the appearance feature is not sufficient to deal with some difficult scenarios, e.g., the same identity dressing different clothing, low resolution videos, the dark illumination cases. Therefore, gait recognition can serve as an effective supplement or alternative to overcome these issues.

There are two popular ways to recognize gaits in literatures, i.e., model-based [1, 2, 3, 4] and appearance-based [5, 6, 7, 8]. The model-based approaches focus on the articulated human features such as the size of a link or joint angles, which can tolerate the appearance changes of an identity due to the clothings or accessories. These approaches require to

preprocess the raw RGB videos to capture the pose structure or silhouettes. On the other hand, several studies have proposed appearance-based gait recognition approaches, which use RGB image sequences as input to recognize the identities directly. However, model-based approaches lose the body shape information and require high accuracy human pose estimation results for gait recognition. Moreover, appearance-based approaches suffer from the sensitivity to the identities' covariates (e.g., dressing and carrying conditions).

In this paper, we propose a novel framework to generate the **Temporal Attention and Keypoint-guided Embeddings** in a principle way called GaitTAKE. The intuition of GaitTAKE is to take both global and local appearance features into account, then the learning of the silhouette embedding is trained by the temporal information. Thus, we can not only solve the flaws by the temporal pooling but also fuse the temporal information into the global and local features. Moreover, we combine the human pose information with the mentioned global and local features so that our method can achieve the large amount of improvement in the wearing coat scenario of gait recognition, which is the most difficult case in gait recognition since the coat will cover most of the area of the human legs. GaitTAKE forms embeddings over multiple frames with a global and local convolutional neural network [9] and human pose information with temporal attention mechanism. According to our experimental results, GaitTAKE achieves the state-of-the-art performance in CASIA-B [10] and OU-MVLP [11] benchmarks.

2. RELATED WORKS

Due to the growth of deep learning, many researchers exploit convolutional neural networks (CNNs) to achieve great improvement for gait recognition [8, 12, 13, 14, 15, 16]. The feature representation ability is robust, e.g., cross-view gait sequence can be recognized only based on the CNN feature and well-designed loss function.

In terms of taking advantage of temporal information, there are two types of deep learning approaches: recurrent neural networks (RNNs) and 3D CNNs. In RNNs, the

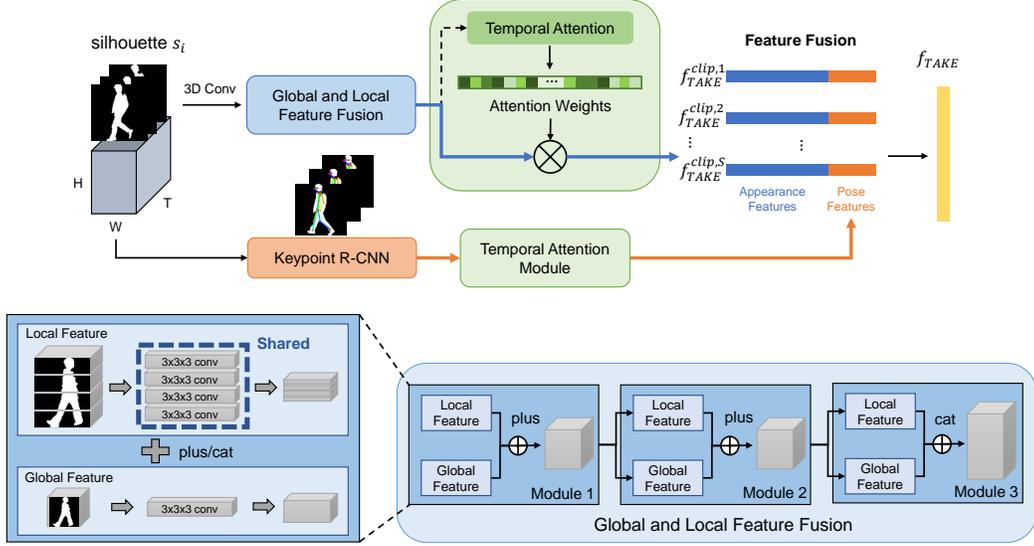


Fig. 1. The architecture of the proposed GaitTAKE framework. First, the frame-level features for the gait videos are extracted by a global and local feature fusion backbone and a human pose feature extraction backbone Keypoint-RCNN. The extracted frame-level features (i.e., appearance features and human pose features) are fed into the temporal attention (TA) module. The TA weights are applied on the frame-level features to obtain the sequence-level features for each gait video.

features are learned by a sequence of consecutive frames [16, 17, 18]. For a 3D CNN, the spatio-temporal information can be extracted by the 3D tensors [15, 19, 9]. Nonetheless, there is a limitation of using 3D CNNs for gait recognition, which is the lack of flexibility for variable-length sequences.

3. PROPOSED METHOD

3.1. TA-based Global and Local Feature Fusion

As shown in Fig. 1, the proposed feature extraction network architecture aims to simultaneously extract the global and local feature information and temporal information from the silhouette images. First, the 3D convolution is applied to extract more representative features from the silhouette images, since 3D convolution network is proved as an effective feature extractor for gait recognition [9, 15, 19]. After that, the local information is extracted by horizontally dividing the image into several body partitions and the global information is extracted by a whole silhouette image. In order to aggregate the temporal information into the local and global feature maps, two 3D convolutions are applied to the local feature maps and the global feature map, separately. The local feature maps (i.e., body partition features) share the same weights of 3D convolutions. According to [9], the generated global and local features can be added into one feature map to ensemble both global and local information. Then, this global and local feature fusion operation is repeated with the same network configuration and different convolution kernels for n times to

generate the more robust global and local fusion features.

The first step of generating the TA-based global and local feature fusion is to first generate the global and local features separately. We use $X \in \mathbb{R}^{c_1 \times T \times h \times w}$ to represent one sequence of silhouette with length T (the image size is $h \times w$), and $\{X_{local}^i | i = 1, \dots, m\}$ denotes the m local gait partition features. c is the channel size of the feature map. Thus, we can express the global gait feature f_{global} as

$$f_{global}(X) = \phi_{global}^{3 \times 3 \times 3}(X) \in \mathbb{R}^{c_2 \times T \times h \times w}, \quad (1)$$

where $\phi_{global}^{3 \times 3 \times 3}$ denotes 3D convolution operation with kernel size $3 \times 3 \times 3$. And for local gait feature f_{local} , similar mechanism is applied with shared 3D convolution kernels,

$$\begin{aligned} f_{local}(X) &= f_{local}(\{X_{local}^i | i = 1, \dots, m\}) \\ &= \phi_{local}^{3 \times 3 \times 3}(X_{local}^1) \oplus \dots \oplus \phi_{local}^{3 \times 3 \times 3}(X_{local}^m) \\ &\in \mathbb{R}^{c_2 \times T \times h \times w}, \end{aligned} \quad (2)$$

where $\phi_{local}^{3 \times 3 \times 3}$ is the shared 3D convolutional layer with kernel size $3 \times 3 \times 3$; \oplus indicates the concatenation operation.

The TA fusion module is composed of two different structures of global and local convolutional (GLConv) layers, i.e., GLConvA and GLConvB. Fig. 1 shows that there are n GLConv layers in this module for generating the global and local information fusion feature f_{GL} ($n = 3$). The last GLConv layer is GLConvB and the rest of other GLConv

layers are GLConvA,

$$\begin{aligned} GLConvA(X) &= f_{global}(X) + f_{local}(X) \\ &\in \mathbb{R}^{c_2 \times T \times h \times w}. \end{aligned} \quad (3)$$

$$\begin{aligned} GLConvB(X) &= f_{global}(X) \oplus f_{local}(X) \\ &\in \mathbb{R}^{c_2 \times T \times 2h \times w}. \end{aligned} \quad (4)$$

Therefore, we can apply flatten operation $\xi(\cdot)$ to get the global and local information fused feature f_{GL} ,

$$\begin{aligned} f_{GL} &= \xi(GLConvB(GLConvA(GLConvA(X)))) \\ &\in \mathbb{R}^{T \times D_{GL}}, \end{aligned} \quad (5)$$

where D_{GL} is the dimension of the f_{GL} .

After obtaining the global and local information fused feature f_{GL} , we can start to apply the TA mechanism to generate the final embedding f_{TGL} . First of all, the sequence of each subject is split into several clips. Assume the clip size is L , $S = \lfloor \frac{T}{L} \rfloor$ is the number of clips, and D indicates the dimension of the clip-level feature.

$$f_{GL}^{clip} = \{f_{GL}^{clip,1}, \dots, f_{GL}^{clip,S}\} \in \mathbb{R}^{S \times L \times D}. \quad (6)$$

Then, there are two convolutional layers used for each clip in the TA module $\mathcal{T}_{GL}(\cdot)$ to produce a feature vector. We subsequently apply a softmax layer to this feature vector to generate a $1 \times L$ -dim attention vector \mathcal{A}_{GL} for weighting the frame-level feature so that the clip-level feature $f_{TGL}^{clip,i} \in \mathbb{R}^{1 \times D}$ can be created.

$$f_{TGL}^{clip,i} = \mathcal{T}_{GL}(f_{GL}^{clip,i}) = \mathcal{A}_{GL} \cdot f_{GL}^{clip,i} \in \mathbb{R}^{1 \times D}. \quad (7)$$

$$\mathcal{A}_{GL} = \sigma_{GL}(\delta_{GL,2}(\delta_{GL,1}(f_{GL}^{clip,i}))) \in \mathbb{R}^{1 \times L}. \quad (8)$$

where $\sigma_{GL}(\cdot)$ is the softmax operation; $\delta_{GL,1}$ and $\delta_{GL,2}$ denote the first and second convolutional layer, respectively.

Finally, one average pooling layer $\psi_{GL}(\cdot)$ is applied to these clip-level embeddings f_{TGL}^{clip} to generate the final embedding f_{TGL} .

$$f_{TGL} = \psi_{GL}(f_{TGL}^{clip}) \in \mathbb{R}^{1 \times D}. \quad (9)$$

3.2. Temporal Aggregated Human Pose Feature

In our framework, we not only consider the appearance embedding feature but also the human pose features since gait recognition is significantly related to the corresponding human pose. We use the keypoint R-CNN [20] to obtain the human pose information. Since not all of the gait recognition datasets contain the human pose information, we use a pretrained model which is trained on COCO dataset to infer the human pose information based on the available RGB images as the ground truth human pose labels. Then, we use the human pose label to train the keypoint R-CNN based on the silhouette images so that we can use the trained keypoint

R-CNN model to infer the human pose information on the silhouette images.

After the human pose is estimated, we use the resulting 2D keypoints (body joints) as the extra features for the gait recognition. The dimension of the human pose features \mathcal{K} for each frame is 17×3 , where 17 is the number of joints, and 3 denotes the 2D joint coordinates (x, y) and corresponding confidence score c . Similar to the appearance features, we also apply the temporal attention technique on the human pose features to aggregate the frame-level features into the clip-based human pose features, and then concatenate the Temporal Aggregated Human Pose Feature with f_{TGL} as the final representation f_{TAKE}^{clip} for gait recognition.

Consequently, we use a Generalized-Mean pooling (GeM) [9] to integrate the spatial information into the feature maps. GeM can effectively generate the more robust representation from the spatial information. Traditionally, researchers fuse the feature from average pooling and max pooling results by a weighted sum, on the other hand, GeM can directly fuse these two different operations to form a feature map, with $p = 1$ being equal to average pooling and $p = \infty$ being equal to max pooling,

$$f_{GeM} = (\psi_{GeM}((f_{TAKE})^p))^{\frac{1}{p}}, \quad (10)$$

where $\psi_{GeM}(\cdot)$ is an average pooling operation.

3.3. Loss Function

The last step of feature extraction is to apply C multiple different fully connected layers to the same f_{GeM} to generate C one-dimensional embedding f . Thus, each subject can be represented by C different embeddings, and all the f of the subject is used to calculate the loss independently. The loss function of our architecture is triplet loss, which is widely used and proved to have superior performance in ReID tasks.

The definition of the triplet loss function is as follows:

$$l_{triplet}(a) = \left[m + \sum_{p \in P(a)} w_p D_{ap} - \sum_{n \in N(a)} w_n D_{an} \right]_+, \quad (11)$$

where m is the margin, D_{ap} and D_{an} indicates the distances between the anchor sample a to form the positive instance and negative instance, respectively. Moreover, w_p and w_n mean the weights of positive and negative instances.

4. EXPERIMENTS

In this work, we use two benchmarks for evaluating the proposed GaitTAKE, namely CASIA-B and OU-MVLP. The first part of this section is to describe the details of the implementation. The second part is to compare GaitTAKE with other state-of-the-art methods in these two datasets.

Setting	Probe	Method	Probe View											Mean
			0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°	
LT(74)	NM#5-6	Gaitset [13]	90.8	97.9	99.4	96.9	93.6	91.7	95.0	97.8	<u>98.9</u>	96.8	85.8	95.0
		MT3D [15]	<u>95.7</u>	98.2	99.0	97.5	95.1	93.9	96.1	<u>98.6</u>	99.2	<u>98.2</u>	<u>92.0</u>	96.7
		GaitGL [9]	94.6	97.3	98.8	97.1	<u>95.8</u>	<u>94.3</u>	<u>96.4</u>	98.5	98.6	<u>98.2</u>	90.8	96.4
		GaitPart [14]	94.1	<u>98.6</u>	<u>99.3</u>	98.5	94.0	92.3	95.9	98.4	99.2	97.8	90.4	96.2
		GaitTAKE (Ours)	96.7	98.6	99.1	<u>98.1</u>	97.3	96.3	98.0	98.9	99.2	99.2	96.4	98.0
	BG#1-2	Gaitset [13]	83.8	91.2	91.8	88.8	83.3	81.0	84.1	90.0	92.2	94.4	79.0	87.2
		MT3D [15]	<u>91.0</u>	<u>95.4</u>	<u>97.5</u>	94.2	<u>92.3</u>	86.9	<u>91.2</u>	<u>95.6</u>	<u>97.3</u>	<u>96.4</u>	86.6	93.0
		GaitGL [9]	90.3	94.7	95.9	94.0	91.9	86.5	90.5	95.5	97.2	96.3	<u>87.1</u>	92.7
		GaitPart [14]	89.1	94.8	96.7	<u>95.1</u>	88.3	<u>94.9</u>	89.0	93.5	96.1	93.8	85.8	91.5
		GaitTAKE (Ours)	96.7	97.0	97.9	97.6	97.9	95.7	97.0	98.2	99.0	99.0	96.4	97.5
	CL#1-2	Gaitset [13]	61.4	75.4	80.7	77.3	72.1	70.1	71.5	73.5	73.5	68.4	50.0	70.4
		MT3D [15]	76.0	87.6	89.8	85.0	81.2	75.7	81.0	84.5	85.4	82.2	68.1	81.5
		GaitGL [9]	<u>76.7</u>	<u>88.3</u>	<u>90.7</u>	<u>86.6</u>	<u>82.7</u>	<u>77.6</u>	<u>83.5</u>	<u>86.5</u>	<u>88.1</u>	<u>83.2</u>	<u>68.7</u>	83.0
		GaitPart [14]	70.7	85.5	86.9	83.3	77.1	72.5	76.9	82.2	83.8	80.2	66.5	78.7
		GaitTAKE (Ours)	89.1	95.3	96.2	93.9	91.5	90.5	92.5	93.3	93.0	91.8	87.0	92.2

Table 1. Rank-1 accuracy (%) of the proposed method on CASIA-B under all views, different size of training data and conditions, excluding identical-view cases. The three walking conditions of sequences include normal (NM), walking with bag (BG) and wearing coat or jacket (CL). The **best** and second accuracy of each probe view will be in bold and underlined respectively.

Method	Probe View														Mean
	0°	15°	30°	45°	60°	75°	90°	180°	195°	210°	225°	240°	255°	270°	
Gaitset [13]	79.5	87.9	89.9	90.2	88.1	88.7	87.8	81.7	86.7	89.0	89.3	87.2	87.8	86.2	87.1
GaitPart [14]	82.6	88.9	90.8	<u>91.0</u>	89.7	89.9	89.5	85.2	<u>88.1</u>	<u>90.0</u>	<u>90.1</u>	<u>89.0</u>	<u>89.1</u>	<u>88.2</u>	88.7
GaitGL [9]	<u>84.3</u>	<u>89.8</u>	<u>90.8</u>	<u>91.0</u>	<u>90.5</u>	<u>90.3</u>	<u>89.9</u>	<u>88.1</u>	87.9	89.6	89.8	88.9	88.9	<u>88.2</u>	89.1
GaitTAKE (Ours)	87.5	91.0	91.5	91.8	91.4	91.1	90.8	90.2	89.7	90.5	90.7	90.3	90.0	89.5	90.4

Table 2. Rank-1 accuracy (%) of the proposed method on OU-MVLP under 14 probe views excluding identical-view cases.

4.1. Implementation Details

In our implementation, the batch size $P \times K$ is set to $8 \times 8 = 64$ in both CASIA-B and OU-MVLP datasets. Following [13], we use 30 frames of each input gait sequence for training and the whole gait sequences are used for extracting gait features in testing. In terms of the number of GLConv layers n , we use 3 GLConv layers (i.e., GLConvA, GLConvA and GLConvB) for CASIA-B dataset. Because OU-MVLP dataset is 20 times larger than the CASIA-B dataset, we use a total of 5 layers, which are 4 GLConvA following by 1 GLConvB layer.

Since the CASIA-B dataset does not contain the keypoint information, we adopt the pre-trained Keypoint R-CNN trained on COCO to estimate the keypoints, which are then used as ground-truth to train the Keypoint R-CNN using masked images, instead of RGB images, as input data. The experimental environment is Python 3.7 and Pytorch 1.7 with one Nvidia GV100.

4.2. Gait Recognition Performance

Evaluation on CASIA-B. We compare our method with state-of-the-art methods: Multiple-Temporal-Scale 3D Convolutional Neural Network (MT3D) [15], Gaitset [13], GaitPart [14] and GaitGL [9] in three different conditions (NM, BG, and CL). The experimental results show that the rank-1

accuracy of the proposed method is higher than GaitGL by about 1.6% and 4.8% in NM and BG, and about 9.2% in CL with the setting of large-scale training (LT, i.e., 74 subjects for training), respectively. It shows that the proposed method has significant advantages in the BG and CL conditions, indicating that the representation of GaitTAKE is much more discriminative than other state-of-the-art methods.

Evaluation on OU-MVLP. We also evaluate the performance of GaitTAKE on the OU-MVLP dataset, where we follow the same training and test protocols as the GaitSet, GaitPart and GaitGL methods for fair comparison. The experimental results are shown in Table 2 and then our method can also achieve the best performance in all cases.

5. CONCLUSION

In this paper, we propose GaitTAKE, which utilizes the temporal attention module to generate the embedding for multi-view gait recognition. We use human pose information and temporal attention to construct the more robust features. Our experimental results show that we can achieve the state-of-the-art performance rank-1 accuracy on two representative gait recognition benchmarks: CASIA-B and OU-MVLP dataset.

6. REFERENCES

- [1] David Kenneth Wagg and Mark S Nixon, "On automated model-based extraction and analysis of gait," in *Sixth IEEE International Conference on Automatic Face and Gesture Recognition, 2004. Proceedings.* IEEE, 2004, pp. 11–16.
- [2] ChewYean Yam, Mark S Nixon, and John N Carter, "Automated person recognition by walking and running via model-based approaches," *Pattern recognition*, vol. 37, no. 5, pp. 1057–1072, 2004.
- [3] David Cunado, Mark S Nixon, and John N Carter, "Automatic extraction and description of human gait models for recognition purposes," *Computer Vision and Image Understanding*, vol. 90, no. 1, pp. 1–41, 2003.
- [4] Gunawan Ariyanto and Mark S Nixon, "Marionette mass-spring model for 3d gait biometrics," in *2012 5th IAPR International Conference on Biometrics (ICB).* IEEE, 2012, pp. 354–359.
- [5] Jinguang Han and Bir Bhanu, "Individual recognition using gait energy image," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 2, pp. 316–322, 2005.
- [6] Dong Xu, Shuicheng Yan, Dacheng Tao, Lei Zhang, Xuelong Li, and Hong-Jiang Zhang, "Human gait recognition with matrix representation," *IEEE Transactions on circuits and systems for video technology*, vol. 16, no. 7, pp. 896–903, 2006.
- [7] Yu Guan, Chang-Tsun Li, and Fabio Roli, "On reducing the effect of covariate factors in gait recognition: a classifier ensemble method," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 7, pp. 1521–1528, 2014.
- [8] Kohei Shiraga, Yasushi Makihara, Daigo Muramatsu, Tomio Echigo, and Yasushi Yagi, "Geinet: View-invariant gait recognition using a convolutional neural network," in *2016 international conference on biometrics (ICB).* IEEE, 2016, pp. 1–8.
- [9] Beibei Lin, Shunli Zhang, Xin Yu, Zedong Chu, and Haikun Zhang, "Learning effective representations from global and local features for cross-view gait recognition," *arXiv preprint arXiv:2011.01461*, 2020.
- [10] Shiqi Yu, Daoliang Tan, and Tieniu Tan, "A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition," in *18th International Conference on Pattern Recognition (ICPR'06).* IEEE, 2006, vol. 4, pp. 441–444.
- [11] Noriko Takemura, Yasushi Makihara, Daigo Muramatsu, Tomio Echigo, and Yasushi Yagi, "Multi-view large population gait dataset and its performance evaluation for cross-view gait recognition," *IPSJ Transactions on Computer Vision and Applications*, vol. 10, no. 1, pp. 4, 2018.
- [12] Zifeng Wu, Yongzhen Huang, Liang Wang, Xiaogang Wang, and Tieniu Tan, "A comprehensive study on cross-view gait based human identification with deep cnns," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 2, pp. 209–226, 2016.
- [13] Hanqing Chao, Yiwei He, Junping Zhang, and Jianfeng Feng, "Gaitset: Regarding gait as a set for cross-view gait recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 8126–8133.
- [14] Chao Fan, Yunjie Peng, Chunshui Cao, Xu Liu, Saihui Hou, Jiannan Chi, Yongzhen Huang, Qing Li, and Zhiqiang He, "Gaitpart: Temporal part-based model for gait recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14225–14233.
- [15] Beibei Lin, Shunli Zhang, and Feng Bao, "Gait recognition with multiple-temporal-scale 3d convolutional neural network," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 3054–3062.
- [16] Xiang Li, Yasushi Makihara, Chi Xu, Yasushi Yagi, Shiqi Yu, and Mingwu Ren, "End-to-end model-based gait recognition," in *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [17] Rijun Liao, Chunshui Cao, Edel B Garcia, Shiqi Yu, and Yongzhen Huang, "Pose-based temporal-spatial network (ptsn) for gait recognition with carrying and clothing variations," in *Chinese Conference on Biometric Recognition*. Springer, 2017, pp. 474–483.
- [18] Dan Liu, Mao Ye, Xudong Li, Feng Zhang, and Lan Lin, "Memory-based gait recognition.," in *BMVC*, 2016, pp. 1–12.
- [19] Weiwei Xing, Ying Li, and Shunli Zhang, "View-invariant gait recognition method by three-dimensional convolutional neural network," *Journal of Electronic Imaging*, vol. 27, no. 1, pp. 013010, 2018.
- [20] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.