

MOTION MATTERS: A NOVEL MOTION MODELING FOR CROSS-VIEW GAIT FEATURE LEARNING

Jingqi Li* Jiaqi Gao* Yuzhen Zhang* Hongming Shan[#] Junping Zhang*[†]

* Shanghai Key Lab of Intelligent Information Processing, School of Computer Science

[#] Institute of Science and Technology for Brain-inspired Intelligence

Fudan University, Shanghai 200433, China

ABSTRACT

As a unique biometric that can be perceived at a distance, gait has broad applications in person authentication, social security and so on. Existing gait recognition methods suffer from changes in viewpoint and clothing and barely consider extracting diverse motion features, a fundamental characteristic in gaits, from gait sequences. This paper proposes a novel motion modeling method to extract the discriminative and robust representation. Specifically, we first extract the motion features from the encoded motion sequences in the shallow layer. Then we continuously enhance the motion feature in deep layers. This motion modeling approach is independent of mainstream work in building network architectures. As a result, one can apply this motion modeling method to any backbone to improve gait recognition performance. In this paper, we combine motion modeling with one commonly used backbone (GaitGL) as GaitGL-M to illustrate motion modeling. Extensive experimental results on two commonly-used cross-view gait datasets demonstrate the superior performance of GaitGL-M over existing state-of-the-art methods.

Index Terms—motion modeling, plug-and-play

1. INTRODUCTION

Silhouette, a standard modality for appearance-based gait recognition, is a binary map generated by segmenting the individual and background. However, the silhouettes among *different individuals* only have subtle variances when the body shapes are similar, inducing the nondiscriminative of the appearance-dependent gait feature. On the contrary, the walking speed and gait cycle are distinguished even though the body shapes look similar among these individuals. Additionally, the silhouettes of *one individual* visually differ when the clothing or viewpoint varies, revealing the vulnerability of the appearance-dependent gait feature. Nevertheless, this person’s motion information, such as speed and gait cycle, remains consistent. Fortunately, this motion information is reflected in the frame-to-frame changes in the sequence of

silhouettes, which can be explored to obtain discriminative and robust gait features.

Recent works mainly aggregate the sequences in different stages. Template-based methods [2–5] compress all silhouettes into one gait template before extracting features, sacrificing the essential temporal information. Set-based methods [6, 7] rather aggregate after the feature extraction stage by pooling. More recently, many new works [1, 8–11] further aggregate the feature sequence in the feature extraction stage using temporal convolution. However, it is hard to extract the motion information through temporal aggregation. As one recent work claimed [12], only relying on temporal convolution is not enough to ensure the uniqueness of the extracted gait feature, let alone the temporal pooling. But, it is noticeable that its theoretical analysis proves that the relationship between adjacent frames can provide the distinguishability of features.

Motivated by these observations, we propose a novel motion modeling for gait recognition, through utilizing the motion information inherent in silhouette sequence and enhancing the motion information in gait representation. Unlike the prior work that employs local self-similarities as the motion information [12], we define the motion information as the holistic temporal changes of all body parts. Our motion modeling method mainly comprises a **Silhouette-level Motion** extractor (**SiMo**), which facilitates silhouette motion encoding, and a **Feature-level Motion** enhancement (**FeMo**), which preserves feature-level motion details. This motion modeling method is applicable to any existing backbone. To better illustrate its usage, we plug the SiMo and FeMo into GaitGL [1], named GaitGL-M. Additionally, the performance of plugging these two modules into GaitSet [6] is presented in experiments (see Table 3).

The contributions of this paper are summarized as follows. 1) We propose a novel motion modeling method to extract the discriminative and robust gait representation. Moreover, this method is independent of network architecture. Thus one can plug it into any existing backbone. 2) We propose two plug-and-play modules in motion modeling, including a silhouette-level motion extractor and feature-

[†]: Corresponding author

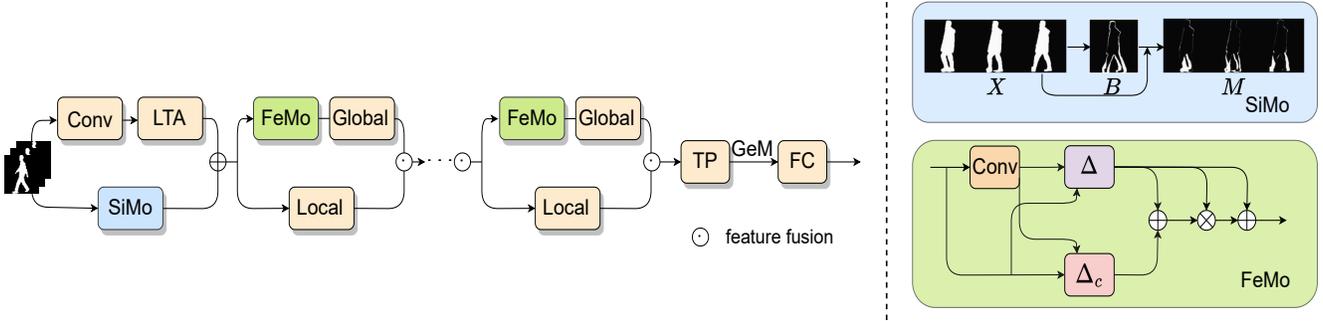


Fig. 1. The picture on the left is the overall framework of GaitGL-M. Removing SiMo and FeMo for motion modeling leaves the architecture of GaitGL [1]. The upper part of the right image is the motion sequence construction process in SiMo, where X is the silhouette sequence, B is the mask, and M represents the generated motion sequence. The lower part is the forward path in FeMo, and the backward direction is treated similarly.

level motion enhancement. Additionally, we combine them with one of the most popular backbones in gait recognition GaitGL as GaitGL-M to show the effectiveness of our proposed modules. 3) Extensive experimental results demonstrate the proposed GaitGL-M’s superiority in the CASIA-B and OU-MVLP datasets, especially when spatial variations appear.

2. METHODOLOGY

Figure 1 presents the framework of GaitGL-M, a combination of our motion modeling and one of the most popular backbones-GaitGL. In GaitGL-M, the silhouette-level motion extractor is paralleled with the original feature extractor, and the features extracted by these two branches are concatenated together as the subsequent input. The feature-level motion enhancement is used to preserve motion details before feature extraction in the global branch.

Before detailing the key modules in GaitGL-M, we provide the description of the gait data. We denote a gait sequence of K frames as $\mathbf{X} = [\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_K] \in \mathbb{R}^{K \times C \times H \times W}$, where C , H , and W represent the number of color channels, height, and width of one frame \mathbf{S} , respectively. Instead of using RGB video frames, the most popular modality used in gait recognition is the silhouette, which is a binary image highlighting the region of the person; correspondingly, $C = 1$ for silhouettes.

2.1. Silhouette-level Motion Extractor (SiMo)

The SiMo explicitly ensures the network’s perception of motion information. It first constructs the motion sequences by seeing the dynamic region exclusively and then extracts the shallow motion feature. Finally, we concatenate the motion feature with the appearance feature. In this way, the sparse motion won’t be diluted by the dense appearance feature.

Motion sequence construction: During a gait cycle, the limbs move alternately. It means that the limbs move for a period of time and are relatively static for a period of time in a gait cycle. Considering the intermittent nature of the limb’s motion, we establish the motion sequence at temporal clips rather than the entire sequence to increase the representational space of motion. In addition, the process of generating silhouette maps contains some noise. Thus, direct performing motion filtering on the whole sequence introduces pseudo-motion information. And selecting the motion signals from temporal clips can also diminish this problem.

More specifically, we uniformly divide a given gait sequence into clips $\mathbf{C} \in \mathbb{R}^{\lfloor K/L \rfloor \times L \times H \times W}$ along the temporal dimension. Then, the motion region mask \mathbf{B}_i for i -th clip can be generated by

$$\mathbf{B}_i = \max(\mathbf{C}_i) - \min(\mathbf{C}_i), \quad (1)$$

where the \max and \min denote the maximum and minimum value of one spatial position among L frames for a clip, and $\mathbf{B}_i \in \mathbb{R}^{H \times W}$. Obviously, we can get the motion clip by multiplying each silhouette in the clip with its corresponding mask in an element-wise manner, defined as follows:

$$\mathbf{M}_i = \mathbf{B}_i \odot \mathbf{C}_i = [\mathbf{B}_i \mathbf{S}_{(i-1)L+1}, \dots, \mathbf{B}_i \mathbf{S}_{(i-1)L+L}]. \quad (2)$$

Then, concatenating the motion clip along the temporal dimension yields the final motion sequence. The constructed motion sequence serves as the supplement input to the original silhouette sequence, which could strengthen the robustness of gait representation against the spatial variants.

Motion feature extraction: In order to extract the shallow motion feature, we first aggregate the frame-level motion within a short clip to obtain $\bar{\mathbf{M}}_i$. Then, all of them are concatenated on the temporal dimension, denoted as $\bar{\mathbf{M}}$. Subsequently, we feed the aggregated motion sequence $\bar{\mathbf{M}}$ into the convolution layer for extracting the motion feature \mathbf{F}^m acquainted with temporal evolution.

2.2. Feature-level Motion Enhancement (FeMo)

The stacked temporal convolutions aggregate the adjacent frames and may weaken the inter-frame differences, making the motion information hard to model. A natural idea to address this issue is to boost the motion-related information before each temporal convolution operation. Inspired by motion encoding in action recognition [13, 14], we regard motion intensity as an attention map to recalibrate the original features and enhance motion-related features. The significant difference from action recognition is that the only action in gait is walking. We aim to mine different walking patterns to identify identities, so fine-grained motion information is necessary.

Firstly, we introduce a bi-directional fine and coarse temporal difference module to distill the subtle motion information from the feature volume. Enlarging the motion search space allows us to understand the importance of different motion directions, and we enlarge the space by spatial convolution.

Consequently, the temporal difference is formulated as:

$$\Delta(\mathbf{G}_t, \mathbf{G}_{t+1}) = \text{Conv}_{2D}(\mathbf{G}_{t+1}) - \mathbf{G}_t, \quad (3)$$

where Conv_{2D} denotes 2D convolution of size 3×3 . Take a single channel for example, the convolution kernel element $w_{i,j}$ semantically represents the importance of different motion directions (i, j) . The larger the value, the higher the degree of attention to motion in the direction. After encoding the fine-grained motion information of each cell, we summarize the spatial information to represent the coarse motion of the whole body by global average pooling:

$$\Delta_C(\mathbf{G}_t, \mathbf{G}_{t+1}) = \text{GAP}(\Delta(\mathbf{G}_t, \mathbf{G}_{t+1})), \quad (4)$$

Moreover, we utilize bi-directional temporal differences to enhance the richness of motion information expression. Overall, the temporal differences are formulated as follows:

$$\begin{cases} \mathbf{D}_t^F = \Delta(\mathbf{G}_t, \mathbf{G}_{t+1}) + \Delta_C(\mathbf{G}_t, \mathbf{G}_{t+1}), \\ \mathbf{D}_t^B = \Delta(\mathbf{G}_{t+1}, \mathbf{G}_t) + \Delta_C(\mathbf{G}_{t+1}, \mathbf{G}_t). \end{cases} \quad (5)$$

Here the superscripts F and B denote the forward and backward operations respectively.

Secondly, we recalibrate the module guided by motion information. A sigmoid function σ is utilized to map the motion intensity into range $(0, 1)$, yielding the average attention from forward and backward directions:

$$\mathbf{W} = (\sigma(\mathbf{D}^F) + \sigma(\mathbf{D}^B))/2 \quad (6)$$

When conducting recalibration, the input feature performs addition with the motion feature, which is the element-wise product of the input feature and motion-aware attention, followed by a convolutional layer to extract motion-aware spatiotemporal feature:

$$\mathbf{G}^m = \text{Conv}(\mathbf{G} + \mathbf{G} \odot \mathbf{W}). \quad (7)$$

2.3. Feature Mapping and Loss Function

After the whole feature extraction stages, we employ temporal max pooling to aggregate the feature volume, followed by generalized-mean pooling (GeM) for spatial pooling [15]. Afterward, we utilize separate fully-connected layer and batch normalization layer to map the feature into a metric space. Following GaitGL [1], we use triplet loss and cross entropy loss function to optimize GaitGL-M.

3. EXPERIMENTS

3.1. Datasets

CASIA-B dataset [16] is a widely used dataset containing 124 individuals. There are 11 camera-perspective uniformly sampling from range $(0^\circ, 180^\circ)$ with 10 sequences in 3 walking conditions for each individual. Normal status (NM) has 6 sequences, bag carrying (BG) and coat-wearing (CL) have 2 sequences respectively. Under the subject-independent protocol [17], we use the large-sample training (LT) strategy [18], in which the sequences are from 74 different identities. **OU-MVLP** is one of the biggest cross-view dataset [19] with 10,307 individuals. There are 14 views sampling from $(0^\circ, 90^\circ)$ and $(270^\circ, 360^\circ)$ respectively per subject and 2 sequences (#seq-00, #seq-01) per view. The train data contains 5,153 individuals, and another 5,154 individuals are taken as test data. In the testing phase, we set #seq-01 as gallery data.

3.2. Implementation Details

The gait silhouettes are normalized before being fed to the network with a fixed input size, 64×44 . And the batch size (p, k) is $(8, 8)$ in CASIA-B dataset and $(32, 8)$ in OU-MVLP dataset, respectively. The optimizer is Adam and the learning rate is $1e-4$ in all experiments. For experiments on CASIA-B, the training iterations is set to 80k and the learning rate decay to $1e-5$ after 70k. For the OU-MVLP, the total iterations are 90k, and the learning rate decay to $1e-5$ after 80k. In our GaitGL-M network for CASIA-B, the number of output channel are 32, 128, 256, 256 for each stage respectively. Since the OU-MVLP is 20 times bigger than CASIA-B, we directly double the convolution layers in each block. Thus the output channel of each stage holds 32, 128, 256, 256. After the first stage, a channel interaction layer implemented by a 1×1 convolution is added on OU-MVLP. Other hyperparameters are following the backbone’s settings. We use four NVIDIA GeForce RTX 3090 GPUs for training GaitGL-M.

3.3. Comparison with the State-of-the-art Method

From Table 1, it can be seen that the average rank-1 accuracies of GaitGL-M outperforms GaitGL by 0.6%, 1.6% and 4.5% in the NM, BG, and CL conditions and implies the superiority of GaitGL-M. Noteworthy, the performance on 90°

Gallery NM#1-4		0°-180°											Mean
Probe		0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°	
NM#5-6	GaitSet [6]	90.8	97.9	99.4	96.9	93.6	91.7	95.0	97.8	98.9	96.8	85.8	95.0
	GaitPart [8]	94.1	98.6	<u>99.3</u>	98.5	94.0	92.3	95.9	98.4	99.2	97.8	90.4	96.2
	GaitEdge [◊] [20]	97.2	99.1	99.2	<u>98.3</u>	97.3	<u>95.5</u>	<u>97.1</u>	99.4	99.3	98.5	96.4	<u>97.9</u>
	GaitGL [1]	96.0	98.3	99.0	97.9	96.9	95.4	97.0	98.9	99.3	<u>98.8</u>	94.0	97.4
	GaitGL-M	<u>96.3</u>	<u>98.8</u>	99.1	98.1	<u>97.2</u>	96.5	98.2	<u>99.1</u>	99.3	99.2	<u>95.9</u>	98.0
BG#1-2	GaitSet [6]	83.8	91.2	91.8	88.8	83.3	81.0	84.1	90.0	92.2	94.4	79.0	87.2
	GaitPart [8]	89.1	94.8	96.7	95.1	88.3	84.9	89.0	93.5	96.1	93.8	85.8	91.5
	GaitEdge [◊] [20]	95.3	97.4	98.4	97.6	<u>94.3</u>	<u>90.6</u>	<u>93.1</u>	97.8	99.1	98.0	95.0	96.1
	GaitGL [1]	92.6	<u>96.6</u>	96.8	95.5	93.5	89.3	92.2	96.5	98.2	96.9	91.5	94.5
	GaitGL-M	<u>93.7</u>	96.4	<u>97.4</u>	<u>97.2</u>	96.2	93.4	95.5	97.8	<u>98.4</u>	<u>97.8</u>	<u>93.1</u>	96.1
CL#1-2	GaitSet [6]	61.4	75.4	80.7	77.3	72.1	70.1	71.5	73.5	73.5	68.4	50.0	70.4
	GaitPart [8]	70.7	85.5	86.9	83.3	77.1	72.5	76.9	82.2	83.8	80.2	66.5	78.7
	GaitEdge [◊] [20]	84.3	<u>92.8</u>	<u>94.3</u>	<u>92.2</u>	<u>84.6</u>	<u>83.0</u>	<u>83.0</u>	<u>87.5</u>	<u>87.4</u>	<u>85.9</u>	<u>75.0</u>	<u>86.4</u>
	GaitGL [1]	76.6	90.0	90.3	87.1	84.5	79.0	<u>84.1</u>	87.0	87.3	84.4	69.5	83.6
	GaitGL-M	<u>79.6</u>	93.4	95.0	92.4	88.4	82.5	86.9	91.4	93.9	90.1	75.3	88.1

Table 1. Averaged rank-1 accuracies on CASIA-B under three different conditions, excluding identical-view cases. The superscript [◊] notes that the input modality is RGB.

Method	Probe view														Mean
	0°	15°	30°	45°	60°	75°	90°	180°	195°	210°	225°	240°	255°	270°	
GaitSet* [6]	78.7	87.4	89.8	90.0	87.8	88.5	97.5	81.3	86.2	88.9	89.1	87.1	87.6	86.1	86.9
GaitPart* [8]	82.1	88.8	90.7	90.8	89.5	89.7	89.1	84.7	87.4	89.9	90.0	88.7	88.9	87.8	88.4
GaitGL* [1]	<u>84.3</u>	<u>89.9</u>	<u>91.1</u>	<u>91.4</u>	<u>90.9</u>	<u>90.6</u>	<u>90.2</u>	<u>88.3</u>	<u>88.1</u>	<u>90.3</u>	<u>90.4</u>	<u>89.6</u>	<u>89.4</u>	<u>88.6</u>	<u>89.5</u>
GaitGL-M	87.1	91.0	91.4	91.8	91.7	91.3	91.1	90.3	89.6	90.7	90.8	90.5	90.2	89.9	90.5

Table 2. Averaged rank-1 accuracies on OU-MVLP, excluding identical-view cases. The superscript * notes the average rank-1 accuracies are reproduced results in our test sets for a fair comparison.

Methods	NM	BG	CL	Mean
GaitGL [1]	97.4	94.5	83.6	91.8
w/ LAGM [12]	96.8	93.1	84.7	91.5
w/ SiMo	98.1	95.8	87.0	93.6
w/ FeMo	97.2	94.7	84.3	92.1
GaitGL-M	98.0	96.1	88.1	94.1
GaitSet [6]	95.0	87.2	70.4	84.2
w/ LAGM [12]	95.0	87.3	73.3	85.2
GaitSet-M	95.6	89.0	73.0	85.9

Table 3. The top half of the table ablates the effect of SiMo and FeMo. The bottom half shows the results of applying the proposed motion modeling to GaitSet.

has been upgraded by our GaitGL-M, exceeding GaitGL by 1.1% (NM), 4.1% (BG) and 3.5% (CL) since the evident legging movement in this viewpoint benefits GaitGL-M. Moreover, although GaitEdge takes the more informative RGB image as input, our GaitGL-M with the silhouette as input still surpasses it, achieving 1.7% higher in the CL. The superior performance indicates that the GaitGL-M has a strong representation ability, even under challenging conditions.

The experimental results on OU-MVLP (Table 2) show great progress than GaitGL by 1.0%, demonstrating the superiority of GaitGL-M. Note that when discarding the illegal sequences, the average rank-1 accuracy will rise to 97.1%.

3.4. Ablation Studies

To explore the contributions of SiMo and FeMo, we design the ablation studies presented in Table 3. It is worth noting

that the proposed GaitGL-M outperforms GaitGL, even only leaving one motion-aware module.

By comparing the designed SiMo and FeMo, we observe that the SiMo contributes more than the FeMo. The possible reason is that the motion information is missed by the smoothed effect of convolution as the network deepens. Therefore, although it boosts the motion-related element, the features themselves contain little motion information.

To further verify the application to other backbones, we apply our SiMo and FeMo to GaitSet, shown on the bottom part of Table 3. It can be found that the averaged rank-1 accuracy has been upgraded by 1.7% with our modules. By comparing with LAGM [12], which regards similarity as motion information, our method based on exploring the pattern of the temporal changes can achieve more performance improvement.

4. CONCLUSION

This paper proposed a novel motion modeling to enjoy the discrimination and robustness of the motion information. Specifically, one silhouette-level motion extractor and feature-level motion enhancement module have been devised to facilitate the motion features in the whole feature extraction stages. Extensive experimental results verify that motion matters in gait recognition and demonstrate the superiority of our motion modeling, which may serve as a plug-and-play module in future model designs.

5. REFERENCES

- [1] Beibei Lin, Shunli Zhang, and Xin Yu, “Gait recognition via effective global-local feature representation and local temporal aggregation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14628–14636.
- [2] Ju Han and Bir Bhanu, “Individual recognition using gait energy image,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 2, pp. 316–322, 2006.
- [3] Chen Wang, Junping Zhang, Jian Pu, Xiaoru Yuan, and Liang Wang, “Chrono-gait image: A novel temporal template for gait recognition,” in *11th European Conference on Computer Vision*, 2010, pp. 257–270.
- [4] Khalid Bashir, Tao Xiang, and Shaogang Gong, “Gait recognition using gait entropy image,” in *3rd International Conference on Imaging for Crime Detection and Prevention*, 2009, pp. 1–6.
- [5] Changhong Chen, Jimin Liang, Heng Zhao, Haihong Hu, and Jie Tian, “Frame difference energy image for gait recognition with incomplete silhouettes,” *Pattern Recognit. Lett.*, vol. 30, no. 11, pp. 977–984, 2009.
- [6] Hanqing Chao, Yiwei He, Junping Zhang, and Jianfeng Feng, “Gaitset: Regarding gait as a set for cross-view gait recognition,” in *The Thirty-Third AAAI Conference on Artificial Intelligence*, 2019, pp. 8126–8133.
- [7] Saihui Hou, Chunshui Cao, Xu Liu, and Yongzhen Huang, “Gait lateral network: Learning discriminative and compact representations for gait recognition,” in *European Conference on Computer Vision*, 2020, pp. 382–398.
- [8] 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. 14213–14221.
- [9] Beibei Lin, Shunli Zhang, and Feng Bao, “Gait recognition with multiple-temporal-scale 3d convolutional neural network,” in *The 28th ACM International Conference on Multimedia*, 2020, pp. 3054–3062.
- [10] Zhen Huang, Dixiu Xue, Xu Shen, Xinmei Tian, Houqiang Li, Jianqiang Huang, and Xian-Sheng Hua, “3D local convolutional neural networks for gait recognition,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14920–14929.
- [11] Xiaohu Huang, Duowang Zhu, Hao Wang, Xinggong Wang, Bo Yang, Botao He, Wenyu Liu, and Bin Feng, “Context-sensitive temporal feature learning for gait recognition,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12889–12898.
- [12] Tianrui Chai, Annan Li, Shaoxiong Zhang, Zilong Li, and Yunhong Wang, “Lagrange motion analysis and view embeddings for improved gait recognition,” in *CVPR*, 2022.
- [13] Boyuan Jiang, Mengmeng Wang, Weihao Gan, Wei Wu, and Junjie Yan, “STM: spatiotemporal and motion encoding for action recognition,” in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, 2019, pp. 2000–2009.
- [14] Yan Li, Bin Ji, Xintian Shi, Jianguo Zhang, Bin Kang, and Limin Wang, “TEA: temporal excitation and aggregation for action recognition,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, 2020, pp. 906–915.
- [15] Filip Radenović, Giorgos Toliás, and Ondřej Chum, “Fine-tuning cnn image retrieval with no human annotation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 7, pp. 1655–1668, 2018.
- [16] 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*, 2006, vol. 4, pp. 441–444.
- [17] Alireza Sepas-Moghaddam and Ali Etemad, “Deep gait recognition: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [18] Hanqing Chao, Kun Wang, Yiwei He, Junping Zhang, and Jianfeng Feng, “Gaitset: Cross-view gait recognition through utilizing gait as a deep set,” *IEEE transactions on pattern analysis and machine intelligence*, 2021.
- [19] 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. 1–14, 2018.
- [20] Junhao Liang, Chao Fan, Saihui Hou, Chuanfu Shen, Yongzhen Huang, and Shiqi Yu, “Gaitedge: Beyond plain end-to-end gait recognition for better practicality,” *arXiv preprint arXiv:2203.03972*, 2022.