

GaitASMS: Gait Recognition by Adaptive Structured Spatial Representation and Multi-Scale Temporal Aggregation

Yan Sun^{1*}, Hu Long¹, Xueling Feng¹, and Mark Nixon²

¹School of Computer Engineering and Science, Shanghai University, 99 Shangda Road, Shanghai, 200444, China

²School of Electronics and Computer Science, University of Southampton, University Road, Southampton, SO17 1BJ, United Kingdom

*Corresponding author(s). E-mail(s): yansun@shu.edu.cn;
Contributing authors: longhu@shu.edu.cn; fengxueling@shu.edu.cn;
msn@ecs.soton.ac.uk;

Abstract

Gait recognition is one of the most promising video-based biometric technologies. The edge of silhouettes and motion are the most informative feature and previous studies have explored them separately and achieved notable results. However, due to occlusions and variations in viewing angles, their gait recognition performance is often affected by the predefined spatial segmentation strategy. Moreover, traditional temporal pooling usually neglects distinctive temporal information in gait. To address the aforementioned issues, we propose a novel gait recognition framework, denoted as GaitASMS, which can effectively extract the adaptive structured spatial representations and naturally aggregate the multi-scale temporal information. The Adaptive Structured Representation Extraction Module (ASRE) separates the edge of silhouettes by using the adaptive edge mask and maximizes the representation in semantic latent space. Moreover, the Multi-Scale Temporal Aggregation Module (MSTA) achieves effective modeling of long-short-range temporal information by temporally aggregated structure. Furthermore, we propose a new data augmentation, denoted random mask, to enrich the sample space of long-term occlusion and enhance the generalization of the model. Extensive experiments conducted on two datasets demonstrate the competitive advantage of proposed method, especially in complex scenes, i.e. BG and CL. On the CASIA-B dataset, GaitASMS achieves the average accuracy of 93.5% and outperforms the baseline on rank-1 accuracies by 3.4% and 6.3%, respectively, in BG and CL. The ablation experiments demonstrate the effectiveness of ASRE and MSTA. The source code is available at <https://github.com/YanSun-github/GaitASMS>.

Keywords: Biometric, Gait recognition, Adaptive Structured Feature, Temporal Aggregation, Deep Learning

1 Introduction

Gait recognition is one of the most popular biometric technologies, which can be employed for human identification at long-distance. Compared with other biometric modalities, *e.g.*, fingerprint, face, and iris, gait does not require cooperation between target subjects and can be hard to disguise. With such advantages, gait recognition has enormous potential in access control, crime investigation, and social security. However, in real-world scenarios, variations like view changes, different wearing conditions, and occlusion [1–3] can cause dramatic changes in gait silhouettes, which are significant challenges to gait recognition.

Nowadays, many deep learning-based gait recognition frameworks have been proposed to generate discriminative gait feature representations [4–7]. **Spatial Feature Extraction: 1)** There

are some studies [8–10] that take the frame-level features as a whole unit for feature extraction. Since the differences between gait silhouettes are tiny, simply using global feature extractors, *e.g.*, global convolution layers and global max/mean pooling are ineffective in capturing fine-grained body information. **2)** The other studies [5–7, 11] are part-based gait recognition methods, which adopt a predefined segmentation strategy to partition the gait silhouettes for fine-grained gait representations. However, these methods only focus on specific body parts, which may limit the ability to capture the global gait features and neglect the relations among different parts. Moreover, since the silhouettes are easily affected by the wearing conditions and the change of viewpoint, it is difficult for the predefined segmentation strategy to effectively focus on the corresponding local features when the contours change dramatically. **Temporal Feature Extraction:** Most state-of-the-art gait recognition methods extract temporal features from gait sequences [5–11]. Gait sequences contain rich temporal information, but these methods only use short-term temporal extractors or simple downsampling pooling operations to extract temporal gait representations. The argument regarding short-term temporal modeling can incorporate the premise that it does not focus on the periodic motion of gait, which is essential but often disrupted by occlusion. Consequently, the use of short-term temporal window functions exhibits a fundamental weakness, which makes these models exhibit poor robustness to occlusion.

To alleviate these issues, we propose a novel gait recognition framework, GaitASMS, which can extract the adaptive spatial representations from global and local features and aggregate long-short-range temporal features in the gait sequences. Specifically, the Adaptive Structured Representation Extraction Module (ASRE) is presented, which adopts an edge-based attention mechanism to extract local fine-grained gait features containing dynamic information and utilizes 3D convolutions for global gait representation extraction. Furthermore, we also propose a temporal aggregate module, called the Multi-Scale Temporal Aggregation Module (MSTA), to achieve long-short-range temporal information aggregation. The captured rich temporal information can effectively compensate for the missing silhouette information caused by occlusion. Transformer-based methods [12] have become increasingly popular in the field of computer vision recently. TransGait [13] employs the transformer module to fuse multi-modal visual information. Compared with the transformer architecture, MSTA is simple but efficient and it does not require massive computation resources for support. Additionally, a new data augmentation method is proposed, denoted as the random mask, which randomly masks some regions of input gait sequences in the target subject level to improve the robustness of the model to occlusion.

The main contributions of this paper can be summarized as the following four aspects:

- 1) **In ASRE**, the local feature extractor adopts an edge-based attention mechanism to obtain adaptive edge masks and extract fine-grained edge gait representations. The global feature extractor is utilized to extract global spatial information as supplementary information to local gait features.
- 2) **In MSTA**, we argue that the short-range temporal feature includes subtle variations in gait and the long-range temporal information contains the robust gait representation for occlusion. Therefore, we introduce the Multi-Scale Temporal Aggregation Module (MSTA) to effectively aggregate the long-short-range temporal information for extracting discriminative temporal features and enhancing the robustness of the model to occlusion.
- 3) **In random mask**, as a novel data augmentation method, it is introduced to enrich the sample space of long-term occlusion and enhance the generalization of the model. Different from traditional masking operations, the random mask exhibits randomness only at the level of target objects, while at the sequence level, the masking regions remain fixed. This design effectively simulates scenarios with long-term occlusion.
- 4) **In GaitASMS**, we adopt adaptive structured spatial representation extraction and multi-scale temporal aggregation to extract distinctive gait features. Extensive experiments conducted on the widely used gait public datasets, CASIA-B and OU-MVLP, demonstrate that GaitASMS outperforms prior SOTA gait recognition methods, especially in occlusion conditions. Abundant ablation experiments conducted on CASIA-B further prove the superiority of the proposed modules.

2 Related Work

Gait Recognition. According to the type of input data, current gait recognition methods can be classified into two types, *i.e.*, model-based and appearance-based. **Model-based** methods [14–18] take the structure of human body as input, like 2D/3D pose data. In theory, model-based gait recognition methods have higher natural robustness to occlusion due to the absence of appearance information. However, considering the difficulty of accurately extracting human pose information from low-resolution images and complex scenes, the practicality of the model-based methods is limited. **Appearance-based** methods [5–7, 11, 19] utilize RGB images or gait silhouettes as input, and attempt to extract discriminative gait features from them. As appearance-based methods do not require extra pose information extraction, it is relatively low-cost and operationally convenient. However, these methods heavily rely on silhouette information, thus they are extremely sensitive to occlusion. Our proposed approach belongs to the appearance-based method and effectively addresses the above issues.

Most state-of-the-art works have taken spatial feature extraction and temporal modeling as the focus. Below, we will provide a detailed introduction to these two forms of feature extraction. **In spatial feature extraction**, recently proposed gait recognition methods can be divided into two categories: global-based and local-based. The global-based methods [19–23] extract gait features from the whole human body. Specifically, GaitSet [19] utilizes 2D CNN to capture frame-level spatial features, and then aggregates spatial gait representations by multiple statistical functions. MT3D [20] uses 3D CNN to obtain spatial-temporal gait features from whole gait sequences. GLN [23] extracts the silhouette-level and set-level features in different stages, and then merges them by the lateral connections in a top-down manner. GaitEdge [24] uses edge detection to pre-process the gait silhouette to capture its shape and texture information, which may be useful in shallow networks, but as the depth of the network increases, it is difficult to truly extract the most discriminatory gait representations without adaptively adjusting the edge segmentation region. The local-based methods [5, 25–27] usually adopt a predefined segmentation strategy to obtain the parts of human silhouettes and then extract fine-grained features from each part. For instance, GaitPart [5] proposes a focal convolution layer, which is used for obtaining fine-grain gait representations. 3DLocal [27] implements a simple but effective form of 3D local CNNs for capturing detailed gait features from multi-scale local parts. GaitStrip [25] learns the local-based spatial features by directly taking each strip of the human body as the basic unit. However, the aforementioned methods all have some issues: **(1)** The global-based methods have difficulty capturing fine-grained spatial features of gait. And the issue becomes more prominent as the network gets deeper. **(2)** The local-based methods encounter challenges in capturing correlation information among local regions and adaptively extracting the most discriminative local features. Particularly, although GaitGL [26] can extract both global and local spatial features, it struggles to adaptively adjust the segmentation strategy according to the degree of occlusion.

In temporal modeling, these approaches can be divided into three categories: GEI-based, Set-based, and Sequence-based. The GEI-based gait recognition methods [21, 28] obtain Gait Energy Image (GEI) by aggregating all temporal information of a sequence. Then, it extracts the final gait representations from the GEI. The set-based methods [19, 22] treat the entire gait sequence as an unordered set to extract the temporal features. Recently, sequence-based gait recognition methods have outperformed alternative approaches, leading to their emergence as the mainstream approach in gait recognition. GaitPart [5] adopts Micro-motion Capture Module (MCM) to capture short-range temporal information from each sequence. GaitGL [26] utilizes Local Temporal Aggregation (LTA) operation to aggregate local temporal information. However, the methods based on local temporal modeling often perform poorly under occlusion conditions because they have difficulty in effectively learning long-term temporal correlation information to compensate for the missing information caused by occlusion.

To address these issues, we propose a novel Adaptive Structured Representation Extraction Module (ASRE) to capture the most dynamic gait patterns and generate adaptive structured spatial representations. And inspired by MG-TCN [29], we first introduce the Multi-Scale Temporal Aggregation Module (MSTA) from Re-ID to gait recognition and effectively achieve the long-short-range temporal modeling. Furthermore, an effective data augmentation approach is proposed to enlarge the sample of occlusion data and enable each module to be trained fully.

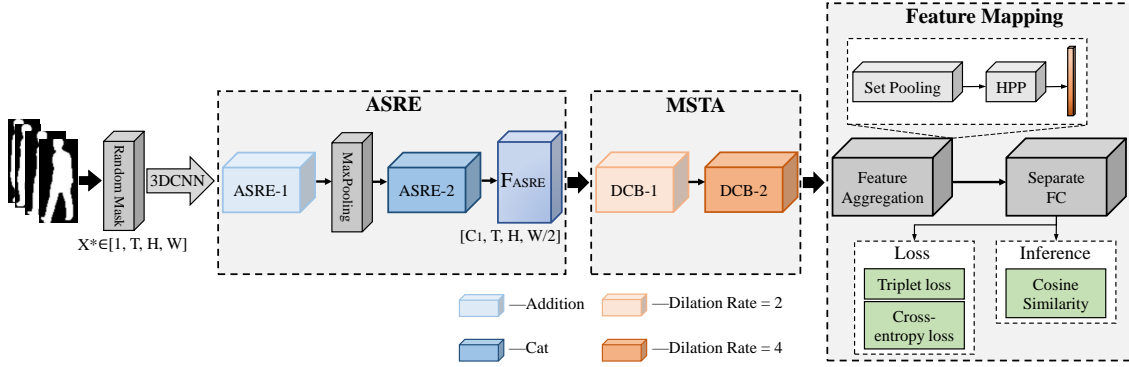


Figure 1: The overview of our GaitASMS. "ASRE" represents the Adaptive Structured Representation Extraction Module. "Addition" means the element-wise addition of local features and global features. "Cat" indicates combining local features and global features in the H dimension. "MSTA" represents the Multi-Scale Temporal Feature Aggregation Module, which is composed of the dilated convolution residual blocks. "HPP" means horizontal pyramid pooling [30].

3 Proposed Method

In this section, we first introduce the pipeline of the proposed method, which can be used to generate discriminative spatial gait representation and extract long-short range temporal correlations. Then the key modules are described, including the Adaptive Structured Representation Extraction Module (ASRE) and the Multi-Scale Temporal Aggregation Module (MSTA). Moreover, a new data augmentation is also proposed. Finally, the details of feature mapping, training, and testing are explained.

3.1 Pipeline

The pipeline of GaitASMS is shown in Fig.1. GaitASMS is composed of two modules, including ASRE and MSTA. Firstly, the sequences of gait silhouettes are fed into the random mask for data augmentation and ASRE extracts adaptive structured spatial representation from the preprocessed data. Max Pooling is added to the framework for obtaining high-level spatial features and reducing the computational cost. Then, to achieve long-short temporal modeling, MSTA extracts the multi-scale temporal features from high-level spatial feature maps. In the end, feature aggregation and separate FC layers are used to map the feature vectors to embedding space for the final individual identification.

Assume that the input of GaitASMS is $X = \{X_1, \dots, X_T\}$, $X \in R^{1 \times T \times H \times W}$, where T is the length of the sequence, H and W are the height and width of each frame, respectively. The pipeline of GaitASMS can be represented as

$$X^* = \text{RandomMask}(X) \quad (1)$$

$$F_{ASRE} = ASRE_2(\text{MaxPooling}(ASRE_1(X^*))) \quad (2)$$

$$F_{MSTA} = DCB_2(DCB_1(F_{ASRE})) \quad (3)$$

$$F_{gait} = FC(FA(F_{MSTA})) \quad (4)$$

where $\text{RandomMask}(\cdot)$ applies masks of the same size to different positions of randomly selected subjects for masking operations. The difference between $ASRE_1(\cdot)$ and $ASRE_2(\cdot)$ is the fusion way of local features and global features. $DCB_1(\cdot)$ and $DCB_2(\cdot)$ represent the dilated convolution blocks, whose dilation rates are 2 and 4, respectively. $FA(\cdot)$ and $FC(\cdot)$ are feature aggregation and feature mapping, respectively. F_{gait} represents the final gait representation.

3.2 Adaptive Structured Spatial Extraction Module

At present, most local-based gait recognition methods only use a predefined segmentation strategy to get the parts of human silhouettes. Although these methods can effectively extract fine-grained features from the parts, they cannot adaptively adjust the size and shape of the local area according to the occlusion status, which undoubtedly limits the representation ability of the network. Thus, the Adaptive Structured Representation Extraction Module (ASRE) is proposed for segmenting the edge of human contours and generating adaptive structured spatial representations. ASRE is shown in Fig.1, which includes the Local Feature Extractor Based on Edge Mask (LEM) and the Global Feature Extractor (GFE). LEM can generate specific edge masks for different sequences, and then the edge masks are employed to obtain local features by masking the whole feature maps. The detailed structure is shown in Fig.2. Thus, ASRE can effectively capture the most distinctive spatial features of each local part, even if under occlusion. GFE is used to extract frame-level spatial information as complementary information of local spatial features.

Inspired by GaitGL [26], we adopt a similar concatenation strategy as GLConv to fuse local spatial features and global spatial features in ASRE. However, the proposed ASRE is largely different from GLFE within GaitGL. Unlike the local feature extractor in GaitGL, which is used to extract details from predefined local maps, the LEM within ASRE generates adaptive structured spatial fine-grained representations through the edge-based attention mechanism. In addition, the subsequent experimental results demonstrate that ASRE can effectively replace GLFE and achieve better performance.

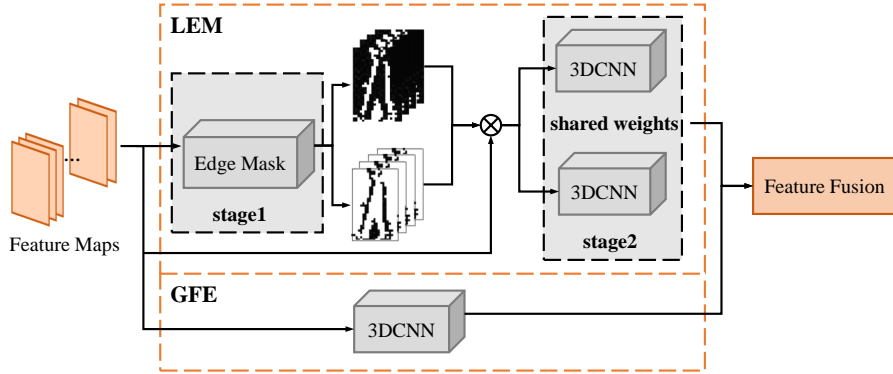


Figure 2: Overview of the ASRE. LEM is the Local Feature Extractor Based on Edge Mask. GFE is the Global Feature Extractor.

Assume that the input feature map is $X^* = \{X_1^*, \dots, X_T^*\}$, $X_i^* \in R^{C \times T \times H \times W}$. ASRE can be represented as

$$F_{ASR} = \begin{cases} add(F_{GFE}, F_{LEM}), & ASRE_1 \\ cat(F_{GFE}, F_{LEM}), & ASRE_2 \end{cases} \quad (5)$$

where F_{GFE} and F_{LEM} are global feature representation and local feature representation, respectively. $cat(\cdot)$ represents the concatenation operation on the H dimension.

As shown in Fig.2, the LEM has two stages. $stage_1$ generates the edge masks for each sequence. Specifically, $MaxPool(\cdot)$ and $AvgPool(\cdot)$ are used as aggregators to extract statistical information in the temporal dimension. Then using $sigmoid(\cdot)$ to normalize the difference between $MaxPool(\cdot)$ and $AvgPool(\cdot)$. The threshold is set as a hyper-parameter to control the intensity of attention on the edge mask. For obtaining richer local feature maps, we also adopt the edge segmentation strategy to generate a complementary mask for the edge mask. $stage_1$ can be formulated as

$$S = sigmoid(MaxPool_{3 \times 1 \times 1}(X^*) - AvgPool_{3 \times 1 \times 1}(X^*)) \quad (6)$$

$$M_{edge} = \begin{cases} 1, & S \geq threshold \\ 0, & S < threshold \end{cases} \quad (7)$$

$$\bar{M}_{edge} = 1 - M_{edge} \quad (8)$$

where \bar{M}_{edge} and M_{edge} is a pair of complementary masks. The generation of M_{edge} is shown in Fig.3. In $stage_2$, a shared weight 3D convolution is employed to extract fine-grained spatial features from the local feature maps based on edge mask. Thus, the F_{LEM} is shown as follows

$$F_{LEM} = 3DConv_{local}^{k \times k \times k} (X^* \otimes M_{edge}) + 3DConv_{local}^{k \times k \times k} (X^* \otimes \bar{M}_{edge}) \quad (9)$$

where $3DConv_{local}^{k \times k \times k}(\cdot)$ is the shared 3D convolution layer with kernel size $k \times k \times k$. \otimes indicates the element-wise multiplication operation. As for the global gait feature F_{GFE} , a similar mechanism is applied with 3D convolution kernels,

$$F_{GFE} = 3DConv_{global}^{k \times k \times k} (X^*) \quad (10)$$

where $3DConv_{global}^{k \times k \times k}(\cdot)$ indicates 3D convolution operation with kernel size $k \times k \times k$.

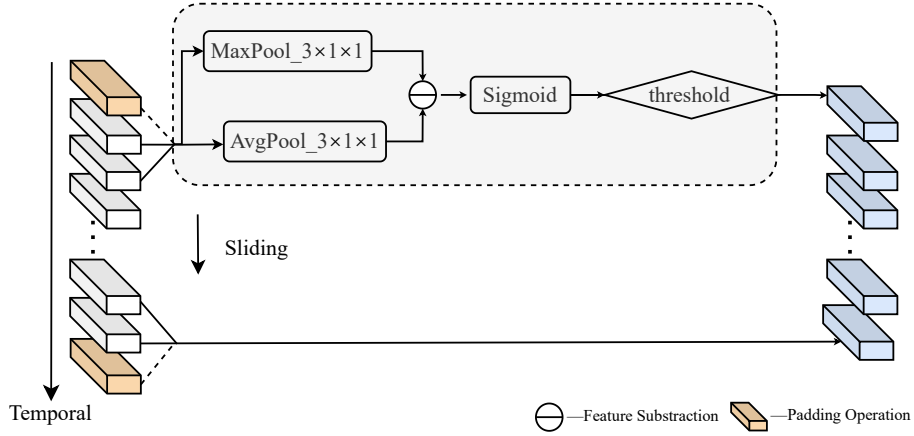


Figure 3: Operation of the Edge Mask.

3.3 Multi-Scale Temporal Aggregation Module

Most of the previous gait recognition methods only use conventional convolution for temporal feature extraction. Due to the limited receptive field, it is difficult to extract long-range temporal relationships. Moreover, the adjacent frames of gait sequences are generally collected in less than 0.04 seconds [31]. In other words, since the difference between adjacent frames is extremely small, it is difficult to learn discriminative temporal features only by temporal modeling of adjacent frames. Thus, we proposed a Multi-Scale Temporal Aggregation Module (MSTA), composed of multi-scale dilated convolution blocks with the residual connection. As shown in Fig.1, the module can efficiently capture long and short-term temporal information. It enables the network to obtain supplemental information through other frames when the contour part is missing, which effectively improves the robustness of the network to occlusion.

The input of MSTA is $F_{ASRE} \in R^{C_1 \times T \times H \times W/2}$, where C_1 represents the channel, $(H, W/2)$ is the size of feature maps. F_{MSTA} can be expressed as

$$F_{MSTA} = DCB_2(DCB_1(F_{ASRE})) \quad (11)$$

where DCB_1 and DCB_2 are dilated convolution blocks, whose dilation rates are 2 and 4, respectively. The detailed pipeline of the DCB is depicted in Fig.4. In each DCB , the input feature maps are processed sequentially through $(Dilated\ Temporal\ 3DConv, Relu, BatchNorm) * 2$, and the residual connection is also employed in DCB .

3.4 Gait Recognition Head

A gait recognition head is employed to map the extracted feature maps into latent embedding space, generating the final gait representation [31]. It generally includes temporal feature mapping and spatial feature mapping.

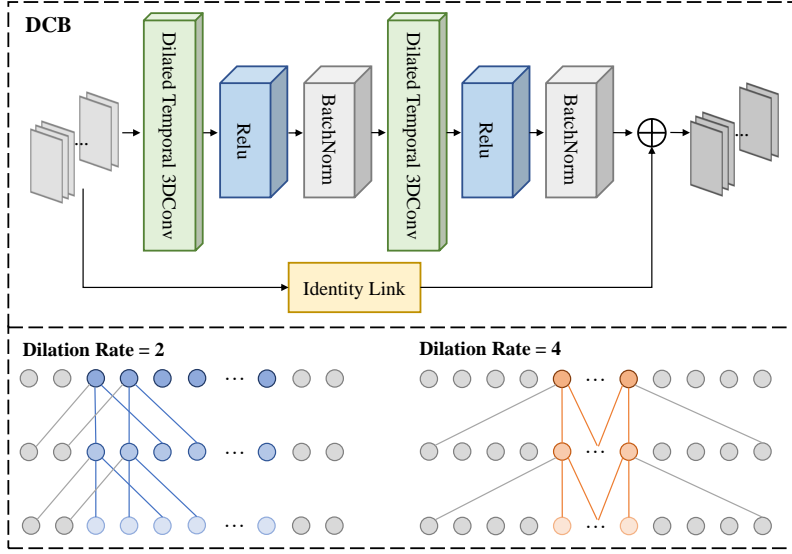


Figure 4: Overview of the DCB. It is the Dilated Convolution Block, which is composed of dilated 3D convolution layers, Relu, and BatchNorm.

For the temporal feature mapping, we introduce a temporal Max-Pooling layer to aggregate the temporal information of variable-length gait sequences [20]. Assuming that the feature map $F_{MSTA} \in R^{C_2 \times T \times H \times W/2}$ is the final output of the feature extraction module. The temporal feature mapping is formulated as

$$Y_{TFM} = MaxPool^{1 \times T \times 1 \times 1}(F_{MSTA}) \quad (12)$$

where $Y_{TFM} \in R^{C_2 \times 1 \times H \times W/2}$ is the output of the temporal feature mapping. $MaxPool^{1 \times T \times 1 \times 1}(\cdot)$ performs a Max-Pooling operation on a sequence of length T .

For the spatial mapping, we use a Generalized-Mean pooling (GeM) [32] to generate multiple horizontal feature representations and then integrate the spatial information into the feature maps. Traditionally, most researchers fuse the features only by a weighted sum of $MaxPool(\cdot)$ and $AvgPool(\cdot)$. However, $GeM(\cdot)$ can directly fuse these two different operations to form a feature map, with $p = \infty$ being equal to $MaxPool(\cdot)$ and $p = 1$ being equal to $AvgPool(\cdot)$,

$$Y_{SFM} = FC(GeM(Y_{TFM})) \quad (13)$$

$$GeM(Y_{TFM}) = \left(AvgPool^{1 \times 1 \times 1 \times W/2}((Y_{TFM})^p) \right)^{\frac{1}{p}} \quad (14)$$

where $AvgPool^{1 \times 1 \times 1 \times W/2}(\cdot)$ is an average pooling operation with kernel size $(1 \times 1 \times 1 \times W/2)$. $FC(\cdot)$ means the Fully Connected layer, which maps gait features into more discriminative embedding space for the final gait recognition.

3.5 Training and Test Details

Training Loss. In this paper, we introduce the combined loss function which consists of the triplet loss [33] and cross-entropy loss to train the proposed network. The triplet loss is utilized to decrease the intra-class distances while increasing the inter-class distance, in addition to the cross-entropy loss employed for classification that facilitates the optimization of the model during training. The combined loss function is represented as

$$L = L_{tri} + L_{cse} \quad (15)$$

where L_{tri} and L_{cse} indicate the triplet loss and cross-entropy loss, separately. L_{tri} is formulated as

$$L_{tri} = \max\{d_p - d_n + margin, 0\} \quad (16)$$

where d_p is the Euclidean distance between positive sample pairs, and d_n is the distance between negative sample pairs. The *margin* is a hyper-parameter that adjusts the optimization difficulty.

Similar to the implementation of OpenGait [34], we use Batch ALL as the sampling strategy. Each batch is formed as $(P \times K)$, where P is the number of subject classes and K denotes the number of samples for a subject class. Due to the memory limitation, the length of gait sequences is set to T during the training stage. The hardware platform consists of an Intel(R) Xeon(R) Silver 4110 CPU @ 2.10GHz, 64GB of RAM, and is equipped with 4 GTX 3090 Ti graphics cards.

Test Stage. In the test stage, we divide the gait datasets into gallery set and probe set for evaluation. It is worth noting that the samples from the gallery set are labeled data, while the samples in the probe set are used for prediction. The proposed network can extract the final gait representation from all frames of the samples belonging to different subject classes. Specifically, we calculate the Euclidean distance between the feature representations in the probe set and the other feature representations in the gallery set. The label of the gallery sample which has the smallest distance from this probe sample can be assigned to the probe sample. Finally, the Rank-1 recognition accuracy is calculated to evaluate the performance.

4 Experiments

In the section, two gait datasets, and the implementation details are first introduced. Then several ablation experiments are performed to show the effectiveness and robustness of GaitASMS.

4.1 Datasets

CASIA-B. Composed of 124 subjects, the CASIA-B [35] is a widely applied gait dataset, in which each subject contains 11 views and each view contains 10 sequences. And the total 10 sequences are obtained under 3 walking conditions, *i.e.*, normal walking (NM), walking with bags (BG), and walking with coats (CL). All video frames under each condition are captured by 11 fixed cameras and recorded with different viewpoints relative to the walking subject. For fairness, this paper strictly follows the popular protocol carried out by [36]. In detail, the first 74 subjects are regarded as train sets, and the remaining 50 subjects are considered as test sets. During the testing phase, only the first four sequences in NM conditions are treated as gallery sets, and the others make up probe sets.

OU-MVLP. The OU-MVLP is one of the largest public gait datasets [37]. However, similar to CASIA-B, it is also a cross-view dataset. It includes 10307 subjects (5153 subjects for training and the rest 5154 subjects for tests). In the dataset, each subject contains 14 views ($0^\circ - 90^\circ; 180^\circ - 270^\circ$) and each view embodies 2 sequences (*#seq-00*, *#seq-01*)[5]. In the testing phase, *#seq-00* is regarded as gallery data.

4.2 Implementation Details

Training details. Common configuration: the gait silhouettes are normalized by [19], and resized to 64×44 . The same as [31], the *margin* in Eqn.15 is set to 0.2 and the parameter p in Eqn.13 is set to 6.5. Adam is taken as the optimizer and the length of each gait sequence T is set to 30. For the CASIA-B: the batch size $(P \times K)$ is (8×8) . The network structure on the CASIA-B is shown in Fig.1. The output channels of "ASRE1", "ASRE2", "DCB1" and "DCB2" are 64, 128, 256, and 256, respectively. The learning rate (λ) is initialized to $1e-4$ for the experiments and the total number of iterations is set to 80k. At 70k iterations, λ will be reduced to $1e-5$. For the OU-MVLP: the batch size $(P \times K)$ is (32×8) . Because of the more subjects in OU-MVLP, we adopt the deeper network setting that doubles the number of the three modules *i.e.*, "ASRE1", "ASRE2" and "DCB1". And the channels of the three modules are 64, 128, and 256. λ is initialized to $4e-4$ and the total number of iterations is set to 130k. At 60k and 110k iterations, λ will be reduced to $4e-5$ and $4e-6$, respectively.

4.3 Performance Comparison

In this section, we compare our method with state-of-the-art ones on CASIA-B and OU-MVLP datasets, and the main results are given in Table 1 and Table 3, respectively.

Performances on CASIA-B. We follow the dataset scales protocol [26] of CASIA-B and test our method in all views and clothing conditions. The experimental results are shown in Table 1. It can be observed that the proposed method has excellent performance with other SOTA methods in all clothing conditions. In the case of NM, BG and CL, our method outperforms the leading method GaitGL [26] by 0.5%, 1.3%, and 3.1%. Even compared with the bimodal gait recognition method TransGait [13], our method still achieves higher accuracy by 0.9% and 0.9% under the BG and CL conditions, respectively. The recognition accuracy of our method in these conditions achieves 97.9%, 95.8%, and 86.7%, respectively. The results show that our method can effectively extract the discriminative gait features, and further improve the accuracy of the model in complex scenarios, specifically in the BG and CL. According to Table 2, the average accuracy of our method is 93.5%, which surpasses the advanced methods i.e., GaitGL, and CSTL [38] by 1.7% and 1.6%, respectively. Several contributing factors include: **1)** In complex scenarios, ASRE focuses on dynamic gait representations adaptively based on contextual information. **2)** Using the pyramid structure for temporal aggregation allows our model to capture richer temporal information and further improve its robustness to occlusions. **3)** The integration of random mask has extended the sample space of occlusion data, thereby boosting the generalization of our model.

Table 1: Averaged Rank-1 accuracy (%) on CASIA-B per probe angle excluding identical-view cases.

Gallery		0° – 180°											Mean
Probe		0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°	
NM	GaitSet[22]	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[5]	94.1	98.6	99.3	98.5	94.0	92.3	95.9	98.4	99.2	97.8	90.4	96.2
	MT3D[20]	95.7	98.2	99.0	97.5	95.1	93.9	96.1	98.6	99.2	98.2	92.0	96.7
	3DLocal[27]	96.0	<u>99.0</u>	<u>99.5</u>	<u>98.9</u>	<u>97.1</u>	94.2	96.3	99.0	98.8	98.5	95.2	97.5
	GaitGL[26]	96.0	98.3	99.0	97.9	96.9	95.4	97.0	98.9	<u>99.3</u>	98.8	94.0	97.4
	CSTL[38]	<u>97.2</u>	<u>99.0</u>	99.2	98.1	96.2	<u>95.5</u>	<u>97.7</u>	98.7	99.2	98.9	96.5	97.8
	TransGait[13]	97.3	99.6	99.7	99.0	<u>97.1</u>	95.4	97.4	<u>99.1</u>	99.6	98.9	95.8	98.1
	GaitStrip*[25]	96.0	98.4	98.8	97.9	96.6	95.3	97.5	98.9	99.1	<u>99.0</u>	<u>96.3</u>	97.6
	Ours	96.4	98.1	98.6	98.0	97.3	96.7	98.6	99.4	99.2	99.6	95.1	<u>97.9</u>
BG	GaitSet[22]	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[5]	89.1	94.8	96.7	95.1	88.3	84.9	89.0	93.5	96.1	93.8	85.8	91.5
	MT3D[20]	91.0	95.4	<u>97.5</u>	94.2	92.3	86.9	91.2	95.6	97.3	96.4	86.6	93.0
	3DLocal[27]	<u>92.9</u>	95.9	97.8	96.2	93.0	87.8	92.7	96.3	97.9	98.0	88.5	94.3
	GaitGL[26]	92.6	<u>96.6</u>	96.8	95.5	93.5	89.3	92.2	96.5	<u>98.2</u>	96.9	91.5	94.5
	CSTL[31]	91.7	96.5	97.0	95.4	90.9	88.0	91.5	95.8	97.0	95.5	90.3	93.6
	TransGait[13]	94.0	97.1	96.5	96.0	93.5	<u>91.5</u>	<u>93.6</u>	95.9	97.2	97.1	<u>91.6</u>	94.9
	GaitStrip*[25]	92.8	<u>96.6</u>	97.2	96.5	95.2	90.5	93.5	<u>97.5</u>	98.3	97.6	91.4	<u>95.2</u>
	Ours	92.7	96.2	97.3	<u>96.4</u>	95.9	93.4	95.6	98.1	98.3	<u>97.7</u>	91.7	95.8
CL	GaitSet[22]	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[5]	70.7	85.5	86.9	83.3	77.1	72.5	76.9	82.2	83.8	80.2	66.5	78.7
	MT3D[20]	76.0	87.6	89.8	85.0	81.2	75.7	81.0	84.5	85.4	82.2	68.1	81.5
	3DLocal[27]	78.2	90.2	92.0	87.1	83.0	76.8	83.1	86.6	86.8	84.1	70.9	83.7
	GaitGL[26]	76.6	90.0	90.3	87.1	84.5	79.0	84.1	87.0	87.3	84.4	69.5	83.6
	CSTL[38]	78.1	89.4	91.6	86.6	82.1	79.9	81.8	86.3	88.7	86.6	<u>75.3</u>	84.2
	TransGait[13]	80.1	89.3	91.0	89.1	84.7	83.3	85.6	87.5	88.2	<u>88.8</u>	76.6	85.8
	GaitStrip*[25]	<u>79.9</u>	92.3	<u>93.4</u>	<u>89.2</u>	<u>86.0</u>	80.0	<u>86.0</u>	<u>88.5</u>	<u>91.7</u>	87.5	73.5	<u>86.2</u>
	Ours	73.8	<u>91.7</u>	93.5	91.5	87.6	<u>82.6</u>	87.3	91.8	92.9	88.9	72.0	86.7

Superscript * is for the work[25] available on arXiv. Underlining represents the sub-optimal results, while bold represents the best results.

Table 2: The overall average Rank-1 accuracy (%) on CASIA-B.

Method	NM	BG	CL	Mean
GaitSet[22]	95.0	87.2	70.4	84.2
GaitPart[5]	96.2	91.5	78.7	88.8
MT3D[20]	96.7	93.0	81.5	90.4
3DLocal[27]	97.5	94.3	83.7	91.8
GaitGL[26]	97.4	94.5	83.6	91.8
CSTL[38]	97.8	93.6	84.2	91.9
TransGait[13]	98.1	94.9	85.8	92.9
GaitStrip*[25]	97.6	<u>95.2</u>	<u>86.2</u>	<u>93.0</u>
Ours	<u>97.9</u>	95.8	86.7	93.5

Heatmap. To better explain the effectiveness of our modules, the heatmaps on CASIA-B are shown in Fig.5. In $b_{1,2}$ and $c_{1,2}$, it can be observed that ASRE with the edge-based attention mechanism can adaptively focus on the most discriminative edge dynamic features, especially on the leg and arm parts. As shown in the red and blue boxes in Fig.5, even in cases of self-occlusion, limited short-term temporal information and the absence of hand contours, MSTA can still enrich the temporal information and obtain the complete heatmaps through the multi-scale temporal aggregation. For instance, in the first five columns of c_2 , it can be observed that even in the absence of arms, GaitASMS also effectively supplements the missing parts to address occlusion due to the presence of MSTA. These show that the proposed modules can indeed promote the learning of adaptive structured representations and the multi-scale aggregation of temporal features.

Performances on OUMVLP. We also compare our methods with leading video-based approaches, including GEINet, GaitSet, GaitPart, GLN, and GaitGL. Table 3 summarizes the performance of these methods on the OU-MVLP dataset. The results indicate that our method outperforms the others, and achieves the highest accuracy in the majority of views, with an average accuracy of 89.9%. It also shows that our method can be effectively applied to large-scale datasets.

Table 3: Rank-1 accuracy (%) on OU-MVLP under 14 probe views, excluding identical identical-view cases.

Method	Probe view														Mean
	0°	15°	30°	45°	60°	75°	90°	180°	195°	210°	225°	240°	255°	270°	
GEINet[21]	11.4	29.4	41.5	45.5	39.5	41.8	38.9	14.9	33.1	43.2	45.6	39.4	40.5	36.3	35.8
GaitSet[22]	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[5]	82.6	88.9	90.8	91.0	89.7	89.9	89.5	85.2	88.1	90.0	90.1	89.0	89.1	88.2	88.7
GLN[23]	83.8	90.0	91.0	91.2	90.3	90.0	89.4	85.3	89.1	90.5	90.6	89.6	89.3	88.5	89.2
GaitGL[26]	84.9	90.2	91.1	91.5	91.1	90.8	90.3	88.5	88.6	90.3	90.4	89.6	89.5	88.8	89.7
Ours	85.6	90.4	91.2	91.6	91.1	90.9	90.4	89.2	89.1	90.4	90.5	89.8	89.5	89.0	89.9

4.4 Ablation Studies

Effectiveness of ASRE. Unlike most part-based gait recognition methods that only use a fixed segmentation strategy to obtain the parts of human silhouettes, we propose an Adaptive Structured Representation Extraction Module (ASRE), which segments the edges of spatial features in latent space and generates adaptive structured spatial representations. To evaluate the module’s effectiveness, several ablation experiments are designed, referred to as Group A. As shown in Table 4, we substituted ASRE with GLConv from GaitGL and MSTA with a standard 3D temporal convolution in A-b, and the average accuracy of A-b is 91.4%. In contrast, the accuracy of A-e (with ASRE) achieved 91.8%, resulting in an increase of 0.4%. Moreover, a comparison between A-d and A-f shows that ASRE improves accuracy by 1.5%, resulting in a performance of 93.5%. These results demonstrate the effectiveness of ASRE. The reasons are as follows: **1)** It utilizes an edge-based attention mechanism to enable the entire model to selectively focus on fine-grained local gait representations, thereby improving the discriminative power of the model; **2)** It also obtains the global spatial information of gait, which complements the local gait representations and enables our model to learn the key relation between parts.

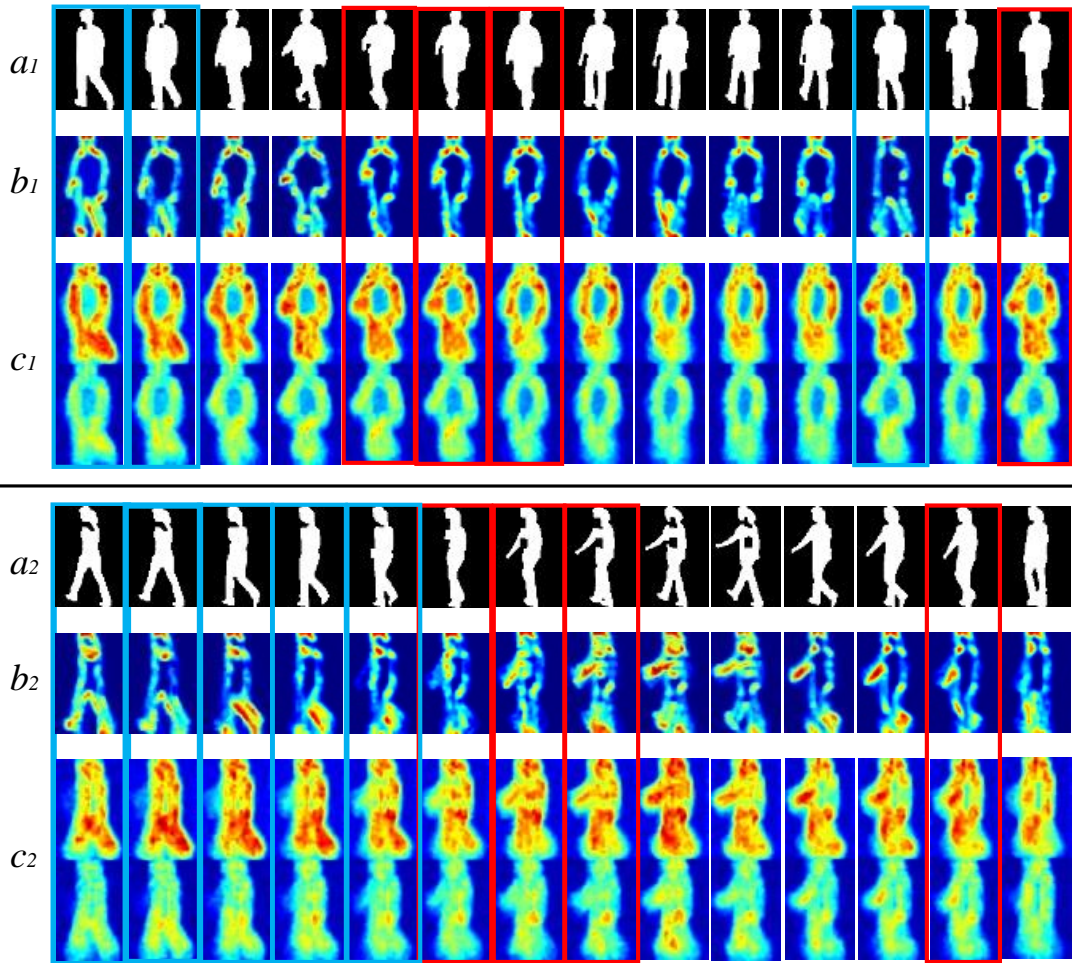


Figure 5: Visualization of the heatmaps for different layers in GaitASMS on CASIA-B. **(a) Top:** the sequence of the silhouettes; **(b) Middle:** the heatmaps of the ASRE-1; **(c) Below:** the heatmaps of the MSTA. The red boxes represent the silhouettes and heatmaps with self-occlusion. The blue boxes represent frames and heatmaps with missing hand contours.

Table 4: Ablation Study: Group A. Control condition: w/ and w/o applying ASRM, MSTA, or Random Mask. Results are rank-1 accuracies (%) on CASIA-B under 11 probe views, excluding identical-view cases.

Group A	Module			Condition			Mean
	ASRM	MSTA	Random Mask	NM	BG	CL	
a	–	–	–	95.9	92.4	80.4	89.6
b	–	–	✓	97.1	94.5	82.7	91.4(↑1.8)
c	✓	–	✓	97.2	94.6	83.7	91.8(↑2.2)
d	–	✓	✓	97.7	94.6	83.7	92.0(↑2.4)
e	✓	✓	–	97.5	95.2	87.2	93.3(↑3.7)
f	✓	✓	✓	97.9	95.8	86.7	93.5(↑3.9)

Table 5: Accuracy (%) of the MSTA with different dilation rates on CASIA-B, under three cloth conditions.

(DCB1, DCB2)	Condition			Mean
	NM	BG	CL	
(1, 2)	97.6	95.1	86.0	92.9
(2, 4)	97.9	95.8	86.7	93.5
(4, 8)	97.9	95.3	86.0	93.1

Effectiveness of MSTA. Most of the previous gait recognition methods [25–27] only use conventional 3D convolution for short-range temporal modeling. However, due to the subtle differences between adjacent frames [31], it is difficult to extract discriminative temporal features solely through the temporal modeling of adjacent frames. Thus, the MSTA module is proposed, composed of multi-scale dilated convolutional blocks with the residual connection. In Table 5, the combination of dilation rates of 2 and 4 in MSTA can better aggregate multi-scale temporal information, as compared to other combinations of dilation rates. As shown in Table 4, compared to A-b, the addition of MSTA in A-d resulted in an average accuracy improvement of 0.6%. In the experiments conducted in A-c (with ASRE) and A-f (with ASRE and MSTA), we observed that the average accuracy of the model also increased by 1.7%. Especially, under BG and CL conditions, the performance improved by 1.2% and 3%, respectively. These results demonstrate that the MSTA module is effective in capturing both long-term and short-term temporal information, thereby enhancing the robustness of the overall model to occlusion.

Effectiveness of the random mask. Unlike traditional data augmentation methods such as horizontal flipping, rotation and erasing, the random mask applies different masks to randomly selected subjects. In Table 6, the accuracy of random erasing shows a decrease of 2.6% compared to the baseline. It suggests that the random erasing struggles to simulate long-term self-occlusion effectively and may result in the absence of gait information. In Table 7, from the ablation experiments of mask rates, the accuracy of a mask rate of 0.1 is the highest. Due to the abundant and easily accessible gait information within the NM condition, its impact is minimal at lower mask rates [0.1, 0.3, 0.5]. At mask rates of [0.7, 0.9], a significant loss of valuable gait information hinders the training of the model and the accuracy has significantly decreased. Therefore, all subsequent experiments are conducted with a mask rate of 0.1. As shown in Table 4, the introduction of the random mask improved the accuracy by 1.8% in the A-a vs. A-b comparison experiments. Especially under the BG and CL conditions, the accuracy improved by 2.1% and 2.3%, respectively. This strongly indicates that the random mask can effectively enhance the robustness of the model to long-term self-occlusion.

Table 6: Accuracy (%) of different data augmentation on CASIA-B, under three cloth conditions.

Data Augmentation	Condition			Mean
	NM	BG	CL	
w/o Baseline	95.9	92.4	80.4	89.6
w/ Random Erasing	94.2	89.8	77.0	87.0
w/ Random Mask	97.1	94.5	82.7	91.4

Table 7: Accuracy (%) of the random mask with different mask rates on CASIA-B, under three cloth conditions.

Mask Rate	Condition			Mean
	NM	BG	CL	
0.1	97.9	95.8	86.7	93.5
0.3	97.8	96.1	86.5	93.5
0.5	97.9	95.8	86.4	93.4
0.7	97.7	95.8	85.3	92.9
0.9	97.3	95.5	84.9	92.6

Generality of the modules. To evaluate the generality of proposed modules, we adapt our

modules to GaitGL. In Table 8, the baseline GaitGL is denoted as B-a, the accuracy of B-b (with ASRE), B-c (with MSTA), and B-d (with random mask) improved by 0.8%, 0.3%, and 0.2%, respectively. Specifically, In B-b and B-c, the GLFE module and LTA module of GaitGL are respectively replaced by ASRE and MSTA. In B-c, the random mask is introduced into GaitGL. These results significantly validate the portability and effectiveness of the proposed modules.

Table 8: Ablation Study: Group B. Control condition: w/ and w/o applying ASRM, MSTA, or Random Mask. This experiment is to verify the generality of ASRM, MSTA, and RM on GaitGL [26].

Group B	Module			Condition			Mean
	ASRM	MSTA	Random Mask	NM	BG	CL	
a	–	–	–	97.4	94.5	83.6	91.8
b	√	–	–	97.6	94.8	85.5	92.6 (↑0.8)
c	–	√	–	97.3	94.4	84.6	92.1 (↑0.3)
d	–	–	√	97.2	94.2	84.6	92.0 (↑0.2)

5 Conclusion

This paper proposes a novel gait recognition framework, denoted GaitASMS, which is based on adaptive structured spatial representation and multi-scale temporal aggregation. By adopting the adaptive edge-based attention mechanism to focus on the most dynamic local edge feature, the GaitASMS can better capture the adaptive structured spatial representations in latent embedding space. Due to the presence of the fundamental weakness caused by the use of short-term temporal window functions, the GaitASMS introduces the multi-scale temporal aggregation module base on dilated 3d convolution to better aggregate the semantic information of gait sequence silhouettes. As a novel data augmentation method, the random mask can enrich the sample space of long-term occlusion and enhance the generalization of the model. Experiments on two well-known public gait datasets, CASIA-B and OU-MVLP, have indicated that compared with other SOTA gait recognition methods, GaitASMS achieves the highest accuracy, especially in occlusion conditions. We hope this work brings the attention of researchers to the adaptive extraction of dynamic local spatial features in gait, as well as the method of utilizing dilated convolutions to achieve both long and short-term temporal aggregation.

Acknowledgments

This work is funded by the National Natural Science Foundation of China, grant number: 62002215; This work is funded by Shanghai Pujiang Program (No. 20PJ1404400).

Declarations

Conflict of interest All authors declare that they have no conflicts of interest.

Data availability The datasets that support the findings of this study are openly available in 10.1109/ICIP.2011.6115889, reference number [35, 37].

References

- [1] A. Sepas-Moghaddam and A. Etemad, “Deep gait recognition: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 264–284, 2021.
- [2] C. Shen, S. Yu, J. Wang, G. Q. Huang, and L. Wang, “A comprehensive survey on deep gait recognition: algorithms, datasets and challenges,” *arXiv preprint arXiv:2206.13732*, 2022.

- [3] Z. Zhang, L. Tran, F. Liu, and X. Liu, “On learning disentangled representations for gait recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 345–360, 2020.
- [4] X. Li, Y. Makihara, C. Xu, Y. Yagi, and M. Ren, “Gait recognition via semi-supervised disentangled representation learning to identity and covariate features,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13309–13319, 2020.
- [5] C. Fan, Y. Peng, C. Cao, X. Liu, S. Hou, J. Chi, Y. Huang, Q. Li, and Z. He, “Gaitpart: Temporal part-based model for gait recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [6] H. Wu, J. Tian, Y. Fu, B. Li, and X. Li, “Condition-aware comparison scheme for gait recognition,” *IEEE Transactions on Image Processing*, vol. 30, pp. 2734–2744, 2021.
- [7] Y. Zhang, Y. Huang, S. Yu, and L. Wang, “Cross-view gait recognition by discriminative feature learning,” *IEEE Transactions on Image Processing*, vol. 29, pp. 1001–1015, 2020.
- [8] T. Wolf, M. Babaee, and G. Rigoll, “Multi-view gait recognition using 3d convolutional neural networks,” in *2016 IEEE International Conference on Image Processing (ICIP)*, pp. 4165–4169, 2016.
- [9] S. Yu, H. Chen, E. B. G. Reyes, and N. Poh, “Gaitgan: Invariant gait feature extraction using generative adversarial networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 532–539, 2017.
- [10] Z. Zhang, L. Tran, X. Yin, Y. Atoum, X. Liu, J. Wan, and N. Wang, “Gait recognition via disentangled representation learning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4710–4719, 2019.
- [11] A. Sepas-Moghaddam and A. Etemad, “View-invariant gait recognition with attentive recurrent learning of partial representations,” *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 3, no. 1, pp. 124–137, 2020.
- [12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [13] G. Li, L. Guo, R. Zhang, J. Qian, and S. Gao, “Transgait: Multimodal-based gait recognition with set transformer,” *Applied Intelligence*, vol. 53, pp. 1–13, 04 2022.
- [14] X. Li, Y. Makihara, C. Xu, Y. Yagi, S. Yu, and M. Ren, “End-to-end model-based gait recognition,” in *Proceedings of the Asian conference on computer vision*, 2020.
- [15] W. An, S. Yu, Y. Makihara, X. Wu, C. Xu, Y. Yu, R. Liao, and Y. Yagi, “Performance evaluation of model-based gait on multi-view very large population database with pose sequences,” *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 2, no. 4, pp. 421–430, 2020.
- [16] W. An, R. Liao, S. Yu, Y. Huang, and P. C. Yuen, “Improving gait recognition with 3d pose estimation,” in *Biometric Recognition: 13th Chinese Conference, CCBR 2018, Urumqi, China, August 11-12, 2018, Proceedings 13*, pp. 137–147, Springer, 2018.
- [17] J. Zheng, X. Liu, W. Liu, L. He, C. Yan, and T. Mei, “Gait recognition in the wild with dense 3d representations and a benchmark,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20228–20237, 2022.
- [18] T. Teepe, A. Khan, J. Gilg, F. Herzog, S. Hörmann, and G. Rigoll, “Gaitgraph: Graph convolutional network for skeleton-based gait recognition,” in *2021 IEEE International Conference on Image Processing (ICIP)*, pp. 2314–2318, 2021.

- [19] H. Chao, Y. He, J. Zhang, and J. Feng, “Gaitset: Regarding gait as a set for cross-view gait recognition,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, pp. 8126–8133, 2019.
- [20] B. Lin, S. Zhang, and F. Bao, “Gait recognition with multiple-temporal-scale 3d convolutional neural network,” in *Proceedings of the 28th ACM international conference on multimedia*, pp. 3054–3062, 2020.
- [21] K. Shiraga, Y. Makihara, D. Muramatsu, T. Echigo, and Y. Yagi, “Geinet: View-invariant gait recognition using a convolutional neural network,” in *2016 international conference on biometrics (ICB)*, pp. 1–8, IEEE, 2016.
- [22] H. Chao, K. Wang, Y. He, J. Zhang, and J. Feng, “Gaitset: Cross-view gait recognition through utilizing gait as a deep set,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 7, pp. 3467–3478, 2021.
- [23] S. Hou, C. Cao, X. Liu, and Y. Huang, “Gait lateral network: Learning discriminative and compact representations for gait recognition,” in *European conference on computer vision*, pp. 382–398, Springer, 2020.
- [24] J. Liang, C. Fan, S. Hou, C. Shen, Y. Huang, and S. Yu, “Gaitedge: Beyond plain end-to-end gait recognition for better practicality,” in *European Conference on Computer Vision*, pp. 375–390, Springer, 2022.
- [25] M. Wang, B. Lin, X. Guo, L. Li, Z. Zhu, J. Sun, S. Zhang, Y. Liu, and X. Yu, “Gaitstrip: Gait recognition via effective strip-based feature representations and multi-level framework,” in *Proceedings of the Asian Conference on Computer Vision*, pp. 536–551, 2022.
- [26] B. Lin, S. Zhang, and X. Yu, “Gait recognition via effective global-local feature representation and local temporal aggregation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14648–14656, 2021.
- [27] Z. Huang, D. Xue, X. Shen, X. Tian, H. Li, J. Huang, and X.-S. Hua, “3d local convolutional neural networks for gait recognition,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14920–14929, 2021.
- [28] C. Zhang, W. Liu, H. Ma, and H. Fu, “Siamese neural network based gait recognition for human identification,” in *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 2832–2836, IEEE, 2016.
- [29] P. Li, P. Pan, P. Liu, M. Xu, and Y. Yang, “Hierarchical temporal modeling with mutual distance matching for video based person re-identification,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 2, pp. 503–511, 2021.
- [30] Y. Fu, Y. Wei, Y. Zhou, H. Shi, G. Huang, X. Wang, Z. Yao, and T. Huang, “Horizontal pyramid matching for person re-identification,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, pp. 8295–8302, 2019.
- [31] B. Lin, S. Zhang, M. Wang, L. Li, and X. Yu, “Gaitgl: Learning discriminative global-local feature representations for gait recognition,” *arXiv preprint arXiv:2208.01380*, 2022.
- [32] F. Radenović, G. Toliás, and O. Chum, “Fine-tuning cnn image retrieval with no human annotation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 7, pp. 1655–1668, 2018.
- [33] A. Hermans, L. Beyer, and B. Leibe, “In defense of the triplet loss for person re-identification,” *arXiv preprint arXiv:1703.07737*, 2017.
- [34] C. Fan, J. Liang, C. Shen, S. Hou, Y. Huang, and S. Yu, “Opengait: Revisiting gait recognition towards better practicality,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9707–9716, 2023.

- [35] S. Zheng, J. Zhang, K. Huang, R. He, and T. Tan, “Robust view transformation model for gait recognition,” in *2011 18th IEEE International Conference on Image Processing*, pp. 2073–2076, 2011.
- [36] Z. Wu, Y. Huang, L. Wang, X. Wang, and T. 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, 2017.
- [37] N. Takemura, Y. Makihara, D. Muramatsu, T. Echigo, and Y. Yagi, “Multi-view large population gait dataset and its performance evaluation for cross-view gait recognition,” *IPSP transactions on Computer Vision and Applications*, vol. 10, pp. 1–14, 2018.
- [38] X. Huang, D. Zhu, H. Wang, X. Wang, B. Yang, B. He, W. Liu, and B. Feng, “Context-sensitive temporal feature learning for gait recognition,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 12909–12918, 2021.