

# MetaGait: Learning to Learn an Omni Sample Adaptive Representation for Gait Recognition

Huanzhang Dou<sup>1</sup>, Pengyi Zhang<sup>1</sup>, Wei Su<sup>1</sup>,  
Yunlong Yu<sup>2\*</sup>[0000-0002-0294-2099], and Xi Li<sup>1,3,4\*</sup>[0000-0003-3023-1662]

<sup>1</sup> College of Computer Science & Technology, Zhejiang University  
{hzdou,pyzhang,weisuzju,xilizju}@zju.edu.cn

<sup>2</sup> College of Information Science & Electronic Engineering, Zhejiang University  
{yuyunlong}@zju.edu.cn

<sup>3</sup> Shanghai Institute for Advanced Study, Zhejiang University

<sup>4</sup> Shanghai AI Laboratory

**Abstract.** Gait recognition, which aims at identifying individuals by their walking patterns, has recently drawn increasing research attention. However, gait recognition still suffers from the conflicts between the limited binary visual clues of the silhouette and numerous covariates with diverse scales, which brings challenges to the model’s adaptiveness. In this paper, we address this conflict by developing a novel MetaGait that learns to learn an omni sample adaptive representation. Towards this goal, MetaGait injects meta-knowledge, which could guide the model to perceive sample-specific properties, into the calibration network of the attention mechanism to improve the adaptiveness from the omni-scale, omni-dimension, and omni-process perspectives. Specifically, we leverage the meta-knowledge across the entire process, where Meta Triple Attention and Meta Temporal Pooling are presented respectively to adaptively capture omni-scale dependency from spatial/channel/temporal dimensions simultaneously and to adaptively aggregate temporal information through integrating the merits of three complementary temporal aggregation methods. Extensive experiments demonstrate the state-of-the-art performance of the proposed MetaGait. On CASIA-B, we achieve rank-1 accuracy of 98.7%, 96.0%, and 89.3% under three conditions, respectively. On OU-MVLP, we achieve rank-1 accuracy of 92.4%.

**Keywords:** Gait recognition, Attention mechanism, Sample adaptive, Learning to learn

## 1 Introduction

As one of the most promising biometric patterns, gait could be recognized at a long distance without the explicit cooperation of humans, thus having wide applications ranging from security check [11], video retrieval [7], to identity identification [3,51]. Most existing approaches [22,45,46] address gait recognition with

---

\* Co-corresponding authors.

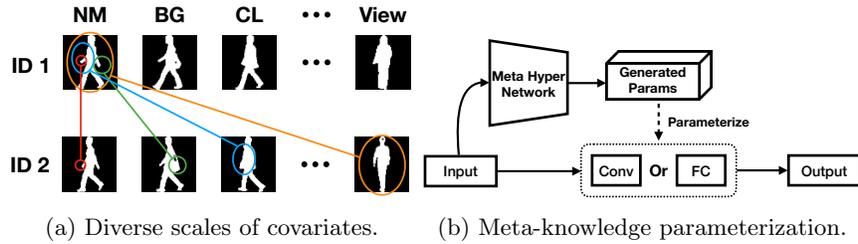


Fig. 1: Illustration of the conflicts and the meta-knowledge. **Left:** The conflicts between limited binary visual clues (colorless and textureless) and numerous covariates with diverse scales, such as bag and clothing, which poses a challenge to the model’s adaptiveness. **Right:** Meta Hyper Network (MHN) learns to learn the meta-knowledge, which could guide the model to perceive sample-specific properties and adaptively parameterize the calibration network.

a two-step process [55]: feature extraction and temporal aggregation. Though significant advances have been achieved, gait recognition still suffers from the conflict between the *limited* binary visual clues (colorless and textureless) and *numerous covariates* with diverse scales of the silhouette shown in Fig. 1a, which poses a huge challenge to the model’s adaptiveness.

Most existing methods tackle this conflict by utilizing the *adaptiveness* of the attention mechanism. For example, the attention mechanism for gait recognition on spatial [36], channel [22], temporal [10], or two of them [35] has been effectively explored. However, the existing attention mechanism still has some limitations, which may harm the adaptiveness. First, the calibration network [34] that performs feature rescaling in the attention mechanism, is static and limited in capturing dependency at a specific scale. Second, the attention mechanism is applied at most two dimensions while leaving one out. Third, only the feature extraction process is considered, while temporal aggregation is ignored.

To address these limitations, we propose a novel framework called MetaGait, to enhance the adaptiveness of the attention mechanism for gait recognition from three perspectives: *omni scale*, *omni dimension*, and *omni process*. The core idea of MetaGait is to leverage *meta-knowledge* [32,68,13], which could guide the model to perceive sample-specific properties, into the calibration network of attention mechanism. Specifically, the meta-knowledge is learned by a Meta Hyper Network (MHN) shown in Fig. 1b and MHN could parameterize the calibration network in a sample adaptive manner instead of being fixed.

Specifically, benefited from the meta-knowledge, we first present Meta Triple Attention (MTA) to adaptively capture the omni-scale dependency in the feature extraction process, leading to the ability to extract walking patterns from diverse scales. The calibration network of MTA is achieved by a weighted dynamic multi-branch structure with diverse receptive fields and parameterized by the meta-knowledge. Second, MTA is designed in homogeneous and applied on spatial/channel/temporal dimensions simultaneously. Third, apart from the fea-

ture extraction process, we present Meta Temporal Pooling (MTP) on temporal aggregation for adaptively integrating temporal information. MTP leverages the meta-knowledge to parameterize an attention-based weighting network, which could excavate the relation between three mainstream temporal aggregation methods with complementary properties (*i.e.*, Max/Average/GeM Pooling [46]). Therefore, MTP could adaptively aggregate their merits for comprehensive and discriminative representation.

Extensive experiments are conducted on two widely used datasets to evaluate the proposed MetaGait framework. The superior results demonstrate that MetaGait outperforms other state-of-the-art methods by a considerable margin, which verifies its effectiveness and adaptiveness.

The major contributions of this work are summarized as follows:

- We present MetaGait framework to address the conflict between limited binary visual clues and numerous covariates with diverse scales. The core idea is to introduce the meta-knowledge learned from Meta Hyper Network to enhance the calibration network’s adaptiveness in the attention mechanism.
- We present Meta Triple Attention (MTA) for the feature extraction process, which aims at adaptively capturing the omni-scale dependency on spatial, channel, and temporal dimensions simultaneously.
- We present attention-based Meta Temporal Pooling (MTP), which could adaptively integrate the merits of three temporal aggregation methods with complementary properties in the temporal aggregation process.

## 2 Related Work

### 2.1 Gait Recognition

**Model-based Approaches.** These methods [2,6,8,26,44,40] aim at modeling the structure of human body from pose information [9,59]. For example, Wang *et al.* [63] propose to use the angle change of body joints to model the walking pattern of different individuals. The advantage of these methods is that they are robust to the clothing and viewpoints conditions. Nevertheless, the model-based approaches suffer from expensive computational costs, accurate pose estimation results, missing ID-related shape information, and extra data collection devices.

**Appearance-based Approaches.** These methods [25,41,52,67,19,27,4,54,42,18,17] learn the features from the silhouette sequences without explicitly modeling the human body structure. For example, GaitSet [10] and GLN [33] deem each silhouette sequence as an unordered set for recognition. GaitPart [22] utilizes 1D convolutions to extract temporal information and aggregate it by a summation or a concatenation. MT3D [45] and 3DLocal [38] propose to exploit 3D convolutions to extract spatial and temporal information at the same time. Appearance-based approaches become popular for their flexibility, conciseness, and effectiveness. The proposed MetaGait is in the scope of appearance-based gait recognition.

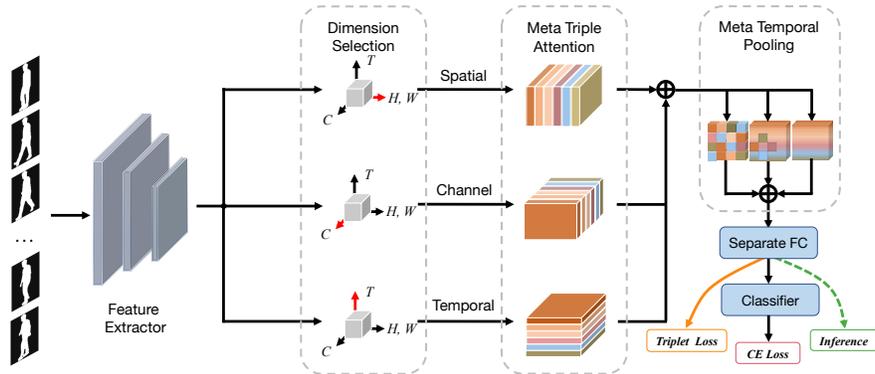


Fig. 2: Overview of MetaGait. Dimension selection refers to the transpose operation for Meta Triple Attention on the selected dimension. Meta Temporal Pooling adaptively aggregates the merits of three temporal aggregation methods with complementary properties. Separate FC is followed by [10,46].

**Attention Mechanism.** Visual attention [34,29,62,65,69], which highlights informative clues and suppresses useless ones, has drawn research attention, and it has been applied to gait recognition successfully. GaitPart [22] performs short-range modeling by channel attention. Zhang *et al.* [70] introduce temporal attention to learn the attention score of each frame by LSTM [20]. Besides, there are methods [36,43] that apply spatial attention. In this paper, we propose to alleviate the conflict between limited visual clues and various covariates with diverse scales from the perspective of the attention mechanism’s adaptiveness.

## 2.2 Dynamic Networks

Dynamic networks can adjust the structures/parameters in an input-dependent manner, leading to several advantages like efficiency, representation power, adaptiveness, and generalizability. Dynamic networks can be mainly divided into dynamic architectures [64,61,49,21,5,39,47] and dynamic parameters [66,12,28,57,14,24,56,15,58]. SkipNet [64] and conv-AIG [61] are two representative approaches to enabling layer skipping to control the architecture. CondConv [66] utilizes the weighted sum of the candidate convolutions according to the input.

Further, Zhang *et al.* [68] point out that dynamic network can be seen as a form of meta-learning [1,32,16,23,68] in learning to learn fashion. In this paper, we leverage the dynamic network for the first time to inject meta-knowledge into the calibration network for improving the model’s adaptiveness.

## 3 Method

In this section, we first present the overview of MetaGait in Fig. 2 and then elaborate on the meta-knowledge learned by Meta Hyper Network (MHN). Further,

we introduce two modules that use meta-knowledge in two separate processes, *i.e.*, Meta Triple Attention for feature extraction and Meta Temporal Pooling for temporal aggregation. Finally, the details of the optimization are described.

### 3.1 Overview

The overview of MetaGait is shown in Fig. 2. First, the gait sequences are fed into the feature extractor, and the feature maps are transposed to perform Meta Triple Attention, which models the omni-scale representation on spatial/channel/temporal dimensions simultaneously. Then, Meta Temporal Pooling adaptively integrates the temporal information with three complementary temporal aggregation methods. Finally, the final objective is computed by the features from separated fully-connected layer [10,46].

### 3.2 Meta Hyper Network

Considering the fact that most attention mechanism in gait recognition applies a static strategy to their calibration network [34], which may harm the model’s adaptiveness, we propose Meta Hyper Network (MHN) to parameterize the calibration network of the attention mechanism adaptively. As shown in Fig. 1b, MHN learns information on data-specific properties of input gait silhouette sequences, *i.e.*, *meta-knowledge*, and generates the parameters of calibration network in a sample adaptive manner.

Given the input  $X \in \mathbb{R}^{C \times T \times H \times W}$ , let  $\mathbf{F}(\cdot)$  be a mapping network, the key to MHN is learning a mapping  $\mathbf{F}_{meta}$  from  $\mathbb{R}^C$  to  $\mathbb{R}^N$  that is used to parameterize the calibration network  $\mathbf{F}_{cali}$  with its parameters  $\mathbf{W}_{cali}$ , *i.e.*, fully connected layer ( $N = C' \times C$ ) or convolution ( $N = C' \times C \times k_h \times k_w \times k_t$ ).  $C$ ,  $C'$ , and  $k$  are the input channel, output channel, and kernel size, respectively. Therefore, the attention mechanism with the meta-knowledge can be formulated as:

$$\mathbf{f} = \mathbf{F}_{cali}(X) \otimes X, \quad s.t. \mathbf{W}_{cali} = \mathbf{F}_{meta}(X), \quad (1)$$

where  $\otimes$  is element-wise multiplication. Specifically, MHN first utilizes Global Average Pooling (GAP) on spatial and temporal dimensions to compute the statics  $m \in \mathbb{R}^{C \times 1 \times 1 \times 1}$  for MHN as:

$$m = GAP(X) = \frac{1}{H \times W \times T} \sum_{i=1}^H \sum_{j=1}^W \sum_{k=1}^T X(i, j, k). \quad (2)$$

Then, the meta-knowledge  $\mathbf{W}_{meta} = \{\mathbf{W}_{meta_1} \in \mathbb{R}^{C \times C}, \mathbf{W}_{meta_2} \in \mathbb{R}^{N \times C}\}$  learned by MHN generates the sample adaptive parameters  $\mathbf{W}_{cali}$  of the calibration network by a Multi-Layer Perceptron (MLP) with Leaky ReLU  $\delta$  as:

$$\mathbf{W}_{cali} = \delta(\mathbf{W}_{meta_2} \delta(\mathbf{W}_{meta_1} m)). \quad (3)$$

In this paper, the meta-knowledge is used to improve the adaptiveness of the attention mechanism on the modules as follows: the global/local calibration stream in Meta Triple Attention (MTA), the soft aggregation gate in MTA, and the weighting network of Meta Temporal Pooling described in Sec. 3.4.

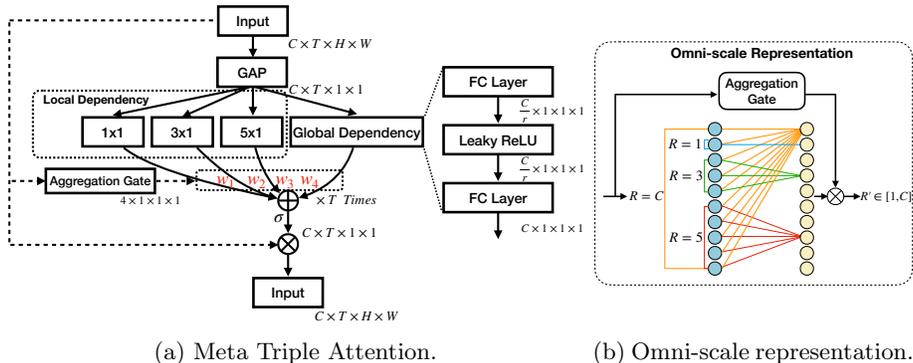


Fig. 3: Illustration of Meta Triple Attention (MTA) and omni-scale representation. (a) MTA is composed of a multi-branch structure with diverse receptive fields weighted by the aggregation gate. (b) MTA achieves omni-scale representation via adaptively weighting the multi-branch structure by soft aggregation gate, leading to the outputs’s receptive field  $R' \in [1, C]$ .

### 3.3 Meta Triple Attention

Though previous attention methods in feature extraction achieve great success, they mainly suffer from two issues. First, they could only capture fixed-scale dependency while numerous covariates present diverse scales, which may harm the model’s adaptiveness. For example, the covariates with small visual changes like bag carrying only require a small receptive field while a large one would bring noises. In contrast, the covariates with significant visual changes like viewpoints require a large receptive field while the small one cannot cover complete visual changes. Second, they only perform attention on two dimensions at most while leaving one dimension out, which is ineffective and insufficient.

To enhance the adaptiveness of the attention mechanism in the feature extraction process, we propose Meta Triple Attention (MTA), which injects the meta-knowledge into its feature rescaling and feature aggregation to capture omni-scale dependency and perform the omni-dimension attention mechanism sufficiently. Thus, MTA differs from the previous attention mechanism in two corresponding aspects: 1) MTA could cope with the numerous covariates at omni scale; 2) MTA performs homogeneous attention mechanism on spatial, channel, and temporal dimensions simultaneously rather than one or two of them. Note that we describe MTA in channel attention on each frame for simplicity while applied in all three dimensions in practice.

Specifically, we leverage the global and local dependency modeling in frame-level with the weighted multi-branch structure to achieve omni-scale representation. As shown in Fig. 3a, for the global channel dependency relation modeling in frame level, a GAP is first applied on the spatial to obtain each frame’s statistics  $s \in \mathbb{R}^{C \times 1 \times 1 \times 1}$ . Then, to effectively capture dimension-wise non-linear global dependency  $\mathbf{f}_{global}$  and evaluate the channel-wise importance, MTA utilizes an

MLP activated by Leaky ReLu, which follows the bottleneck design [34,30] with a dimension reduction ratio  $r$ :

$$\mathbf{f}_{global} = \mathbf{F}_{global}(s) = \mathbf{W}_{g_2} \delta(\mathbf{W}_{g_1} s), \quad (4)$$

where the parameters  $\mathbf{W}_g = \{\mathbf{W}_{g_1} \in \mathbb{R}^{\frac{C}{r} \times C}, \mathbf{W}_{g_2} \in \mathbb{R}^{C \times \frac{C}{r}}\}$  of global calibration stream  $\mathbf{F}_{cali}^{global}$  is adaptively parameterized by MHN.

For local dependency modeling of the calibration network, we design a multi-branch convolutional structure with diverse receptive fields (*i.e.*, kernel size). Therefore, each local calibration stream could capture dependency at a specific scale. To learn omni-scale representation, we propose to aggregate the output  $\mathbf{f}_{global}$  and  $\mathbf{f}_{local}$  of global and local streams in a sample adaptive manner as Eq. (5) instead of being fixed, which is achieved by a soft aggregation gate  $\mathbf{G}$  with meta-knowledge. Next, Sigmoid  $\sigma$  is applied to mapping the values of the attention vector into  $[0, 1]$ :

$$\mathbf{f}_{mta} = \sigma(\mathbf{G}(s)[L+1] * \mathbf{f}_{global} + \sum_{l=1}^L \mathbf{G}(s)[l] * \mathbf{f}_{local}^l) \otimes X, \quad s.t. \ L \geq 1, \quad (5)$$

where  $L$  denotes the number of the receptive field sizes in local stream of the calibration network. The output of the soft aggregation gate  $\mathbf{G}$  is a vector with  $\mathbb{R}^{L+1}$  to weight each stream according to the input. Specifically,  $\mathbf{G}$  is implemented by GAP, an MLP mapping from  $\mathbb{R}^{C \times 1 \times 1 \times 1}$  to  $\mathbb{R}^{(L+1) \times 1 \times 1 \times 1}$ , and Sigmoid in sequential. Therefore, the output's receptive field  $\mathbf{R}'$  is adaptively ranging from 1 to global receptive field  $C$ , which could capture omni-scale dependency.

Besides, previous approaches design heterogeneous attention modules for each dimension to fit the dependency scale that each dimension needs to model. Benefited from the omni-scale dependency modeling, MTA can be efficiently performed on three different dimensions in homogeneous.

### 3.4 Meta Temporal Pooling

To achieve omni-process sample adaptive representation, the meta-knowledge is injected into the temporal aggregation apart from feature extraction. In the recent gait literature [10,22,46], Global Max Pooling (GMP), Global Average Pooling (GAP), and GeM Pooling [53] along the temporal dimension are the mainstream temporal aggregation methods, which represent the salient information, overall information, and an intermediate form between the former two methods, respectively. They can be formulated as:

$$\mathbf{Max}(\cdot) = \mathit{Pool}_{Max}^{T \times 1 \times 1}(\cdot), \quad (6a)$$

$$\mathbf{Mean}(\cdot) = \mathit{Pool}_{Avg}^{T \times 1 \times 1}(\cdot), \quad (6b)$$

$$\mathbf{GeM}(\cdot) = (\mathbf{Mean}(\cdot)^p)^{\frac{1}{p}}, \quad (6c)$$

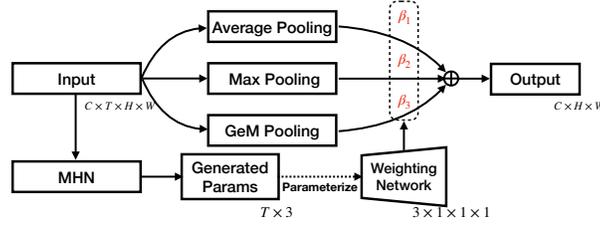


Fig. 4: The illustration of Meta Temporal Pooling (MTP), which aims at leveraging meta-knowledge from MHN to adaptively integrate the merits from three complementary temporal aggregation methods.

where  $p$  in GeM Pooling is a learnable parameter. Though these temporal aggregation methods have been validated their effectiveness individually, their relation is under-explored. We argue that different temporal aggregation methods have their own merits and complementarities to each other, which can be excavated to adaptively integrated temporal clues according to the properties of inputs. Specifically, GMP preserves the most salient information along the temporal dimension while ignoring the majority of information. By contrast, GAP includes the overall temporal information, but the salient information would be diluted out. Though GeM, an intermediate form, can obtain salient temporal information while preserving overall one, it is less robust than GAP and GMP due to the learning stability of unconstrained learnable parameter  $p$ .

To fully exploit their merits, we leverage the meta-knowledge learned by MHN to adaptively integrate the features produced by three temporal aggregation methods as shown in Fig. 4. In detail, we first compute the statics using GAP on spatial and channel dimensions, and we utilize MHN to generate the parameter ( $N = 3T$ ), which is used to parameterize the weighting network, *i.e.*, an FC layer with  $\mathbf{W}_t \in \mathbb{R}^{3 \times T}$  followed by a Sigmoid. Therefore, the weights of different temporal aggregation methods  $\beta \in \mathbb{R}^3$  can be obtained as:

$$\beta = \sigma(\mathbf{W}_t(\text{GAP}(\mathbf{f}_{mta}))). \quad (7)$$

Then,  $\beta$  adaptively weights the features of three complementary temporal aggregation methods and obtain omni sample adaptive representation  $\mathbf{f}_{omni}$  as:

$$\mathbf{f}_{omni} = \beta_1 \text{Mean}(\mathbf{f}_{mta}) + \beta_2 \text{Max}(\mathbf{f}_{mta}) + \beta_3 \text{GeM}(\mathbf{f}_{mta}) \quad (8)$$

### 3.5 Optimization

Following the optimization strategy[46,33,35], we apply Triplet Loss [31]  $\mathcal{L}_{tri}$  and Cross-Entropy loss  $\mathcal{L}_{ce}$  on each horizontal feature independently to train our model as Eq. (9). The similarity metric is set to Euclidean distance.

$$\mathcal{L}_{total} = \mathcal{L}_{tri} + \mathcal{L}_{ce}. \quad (9)$$

## 4 Experiments

### 4.1 Datasets

**CASIA-B** [67]. It is composed of 124 IDs, each of which has 10 groups of sequences, *i.e.*, 6 normal walking (NM), 2 walking with a bag (BG), 2 walking in coats (CL). The views are uniformly distributed in  $[0^\circ, 180^\circ]$ . For evaluation, the protocol is adopted as [10], *i.e.*, small-scale training (ST), medium-scale training (MT), and large-scale training (LT). These three settings select the first 24/62/74 IDs as the training set and the rest 100/62/50 IDs as the test set, respectively. During the evaluation, the first four sequences of each ID under NM are deemed as the gallery, and the rest are used as the probe.

**OU-MVLP** [60]. It is the largest dataset consisting of 10,307 IDs. In OU-MVLP, there are 1 waling condition (NM) with 2 sequences and 14 views, which are uniformly distributed between  $[0^\circ, 90^\circ]$  and  $[180^\circ, 270^\circ]$ . The training set and test set are composed of 5,153 IDs and 5,154 IDs, respectively. For evaluation, the first sequence of each ID is adopted as the gallery, and the rest is the probe.

### 4.2 Implementation Details

**Hyper-parameters.** 1) The resolution of the silhouette is resized to  $64 \times 44$  or  $128 \times 88$  following [33,38,35]; 2) In a mini-batch, the number of the IDs and the sequences of each ID is set to (8, 8) for CASIA-B and (32, 8) for OU-MVLP; 3) Adam optimizer is used with a learning rate of  $1e-4$ ; 4) We train our model for 100k iterations for CASIA-B and 250k for OU-MVLP, where the learning rate is reduced to  $1e-5$  at 150k iterations; 5) The margin of Triplet loss is set to 0.2; 6) The reduction ratio  $r$  in this paper is all set to 2.

**Training Details.** 1) The feature extractor is following the global and local backbone in [46]; 2) The channels of the feature extractor in the three stages are set to (32, 64, 128) for CASIA-B and double for OU-MVLP. 3) The local stream of MTA is implemented with Conv1d and Conv2d for channel/temporal and spatial dimensions, respectively. The receptive fields of the local stream are set to {1,3,5}. Refer to supplementary materials for more details.

### 4.3 Comparison with State-of-the-Art Methods

**Results on CASIA-B.** To evaluate MetaGait on cross-view and large resolution scenarios, we conduct a comparison between MetaGait and latest SOTA as shown in Tab. 1, where MetaGait outperforms SOTA at most views and both two resolutions. Specifically, under NM/BG/CL conditions, MetaGait outperforms previous methods by **0.3%/0.4%**, **0.7%/0.5%**, and **2.7%/2.3%** at the resolution of  $64 \times 44/128 \times 88$  **at least**. Further, MetaGait achieves rank-1 accuracies over **98%** and **96%** under NM and BG, respectively. More importantly, the considerable performance gain on the most challenging condition CL narrows the gap between the performance of NM and CL to less than **10%**, which verifies the robustness of MetaGait under the cross-walking-condition scenario.

Table 1: Averaged rank-1 accuracy on CASIA-B, excluding identical views cases.

Gallery NM #1-4			0°-180°											Mean		
Prob.	Res.	Method	0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°			
NM	64 × 44	GaitSet	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	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	95.7	98.2	99.0	97.5	95.1	93.9	96.1	98.6	99.2	98.2	92.0	96.7		
		CSTL	97.2	99.0	99.2	98.1	96.2	95.5	<b>97.7</b>	98.7	99.2	98.9	96.5	97.8		
		3DLocal	96.0	99.0	99.5	98.9	97.1	94.2	96.3	99.0	98.8	98.5	95.2	97.5		
		GaitGL	96.0	98.3	99.0	97.9	96.9	95.4	97.0	98.9	99.3	98.8	94.0	97.4		
		<b>MetaGait</b>	<b>97.3</b>	<b>99.2</b>	<b>99.5</b>	<b>99.1</b>	<b>97.2</b>	<b>95.5</b>	97.6	<b>99.1</b>	<b>99.3</b>	<b>99.1</b>	<b>96.7</b>	<b>98.1</b>		
	128 × 88	GaitSet	91.4	98.5	98.8	97.2	94.8	92.9	95.4	97.9	98.8	96.5	89.1	95.6		
		GLN	93.2	99.3	99.5	98.7	96.1	95.6	97.2	98.1	99.3	98.6	90.1	96.9		
		CSTL	97.8	99.4	99.2	98.4	97.3	95.2	96.7	98.9	99.4	99.3	96.7	98.0		
		3DLocal	97.8	99.4	99.7	99.3	97.5	96.0	98.3	99.1	99.9	99.2	94.6	98.3		
		<b>MetaGait</b>	<b>98.1</b>	<b>99.4</b>	<b>99.8</b>	<b>99.4</b>	<b>97.6</b>	<b>96.7</b>	<b>98.5</b>	<b>99.3</b>	<b>99.9</b>	<b>99.6</b>	<b>97.0</b>	<b>98.7</b>		
		BG	64 × 44	GaitSet	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	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	91.0			95.4	97.5	94.2	92.3	86.9	91.2	95.6	97.3	96.4	86.6	93.0		
CSTL	91.7			96.5	97.0	95.4	90.9	88.0	91.5	95.8	97.0	95.5	90.3	93.6		
3DLocal	92.9			95.9	<b>97.8</b>	96.2	93.0	87.8	92.7	96.3	97.9	98.0	88.5	94.3		
GaitGL	92.6			96.6	96.8	95.5	93.5	89.3	92.2	96.5	98.2	96.9	91.5	94.5		
<b>MetaGait</b>	<b>92.9</b>			<b>96.7</b>	97.1	<b>96.4</b>	<b>94.7</b>	<b>90.4</b>	<b>92.9</b>	<b>97.2</b>	<b>98.5</b>	<b>98.1</b>	<b>92.3</b>	<b>95.2</b>		
128 × 88	GaitSet		89.0	95.3	95.6	94.0	89.7	86.7	89.7	94.3	95.4	92.7	84.4	91.5		
	GLN		91.1	97.7	97.8	95.2	92.5	91.2	92.4	96.0	97.5	95.0	88.1	94.0		
	CSTL		95.0	96.8	97.9	96.0	94.0	90.5	92.5	96.8	97.9	<b>99.0</b>	<b>94.3</b>	95.4		
	3DLocal		94.7	98.7	98.8	97.5	93.3	91.7	92.8	96.5	98.1	97.3	90.7	95.5		
	<b>MetaGait</b>		<b>95.1</b>	<b>98.9</b>	<b>99.0</b>	<b>97.8</b>	<b>94.0</b>	<b>92.0</b>	<b>92.9</b>	<b>96.9</b>	<b>98.2</b>	98.4	93.5	<b>96.0</b>		
	CL		64 × 44	GaitSet	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	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		76.0		87.6	89.8	85.0	81.2	75.7	81.0	84.5	85.4	82.2	68.1	81.5		
CSTL		78.1		89.4	91.6	86.6	82.1	79.9	81.8	86.3	88.7	86.6	75.3	84.2		
3DLocal		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		76.6		90.0	90.3	87.1	84.5	79.0	84.1	87.0	87.3	84.4	69.5	83.6		
<b>MetaGait</b>		<b>80.0</b>		<b>91.8</b>	<b>93.0</b>	<b>87.8</b>	<b>86.5</b>	<b>82.9</b>	<b>85.2</b>	<b>90.0</b>	<b>90.8</b>	<b>89.3</b>	<b>78.4</b>	<b>86.9</b>		
128 × 88		GaitSet	66.3	79.4	84.5	80.7	74.6	73.2	74.1	80.3	79.7	72.3	62.9	75.3		
		GLN	70.6	82.4	85.2	82.7	79.2	76.4	76.2	78.9	77.9	78.7	64.3	77.5		
		CSTL	84.1	92.1	91.8	87.2	84.4	81.5	84.5	88.4	91.6	91.2	79.9	87.0		
		3DLocal	78.5	88.9	91.0	89.2	83.7	80.5	83.2	84.3	87.9	87.1	74.7	84.5		
		<b>MetaGait</b>	<b>87.8</b>	<b>94.6</b>	<b>93.5</b>	<b>90.3</b>	<b>87.1</b>	<b>84.3</b>	<b>86.1</b>	<b>89.7</b>	<b>93.9</b>	<b>93.4</b>	<b>81.7</b>	<b>89.3</b>		

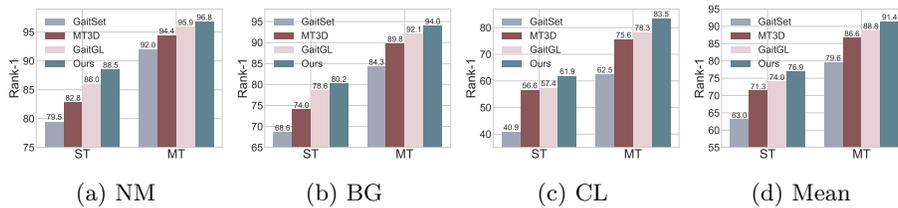


Fig. 5: Comparison with state-of-the-art methods under ST/MT setting.

Then, we evaluate MetaGait under the data-limited scenarios following the protocol in [10]. As the experimental results are shown in Fig. 5, MetaGait outperforms state-of-the-art methods with a significant margin, which further shows the efficiency and robustness of MetaGait under small data scenarios.

Table 2: Comparison with SOTA methods of rank-1 accuracy (%) and mAP (%).

Method	Pub.	Rank-1				mAP
		NM	BG	CL	Mean	
GaitSet [10]	AAAI19	95.0	87.2	70.4	84.2	86.2
GaitPart [22]	CVPR20	96.2	91.5	78.7	88.8	88.7
GLN [33]	ECCV20	96.9	94.0	77.5	89.5	89.2
MT3D [48]	ACM MM20	96.7	93.0	81.5	90.4	90.1
CSTL [37]	ICCV21	97.8	93.6	84.2	91.9	-
3DLocal [38]	ICCV21	97.5	94.3	83.7	91.8	-
GaitGL [46]	ICCV21	97.4	94.5	83.6	91.8	91.5
<b>MetaGait</b>	-	<b>98.1</b>	<b>95.2</b>	<b>86.9</b>	<b>93.4</b>	<b>93.2</b>

Table 3: Averaged rank-1 accuracy on OU-MVLP across different views excluding identical-view cases.

Method	Probe View													Mean	
	0°	15°	30°	45°	60°	75°	90°	180°	195°	210°	225°	240°	255°		270°
GEINet	11.4	29.1	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	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	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	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
CSTL	87.1	91.0	91.5	91.8	90.6	90.8	90.6	89.4	90.2	90.5	90.7	89.8	90.0	89.4	90.2
3DLocal	86.1	91.2	92.6	92.9	92.2	91.3	91.1	86.9	90.8	92.2	92.3	91.3	91.1	90.2	90.9
GaitGL	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
MetaGait (64 × 44)	88.2	92.3	93.0	93.5	93.1	92.7	92.6	89.3	91.2	92.0	92.6	92.3	91.9	91.1	91.9
<b>MetaGait (128 × 88)</b>	<b>88.5</b>	<b>92.6</b>	<b>93.4</b>	<b>93.7</b>	<b>93.8</b>	<b>93.0</b>	<b>93.3</b>	<b>90.1</b>	<b>91.7</b>	<b>92.4</b>	<b>93.3</b>	<b>92.9</b>	<b>92.6</b>	<b>91.6</b>	<b>92.4</b>

Further, to evaluate the comprehensive retrieval performance of MetaGait, we present the average rank-1/mAP performance in Tab. 2, where mAP is computed by the reproduced methods. Specifically, MetaGait outperforms GaitSet by **9.2%/7%**, GaitPart by **4.6%/4.5%**, and GaitGL by **1.6%/1.7%**, which indicates the superior retrieval performance of MetaGait.

**Results on OU-MVLP.** To verify the effectiveness of MetaGait on the large dataset, we evaluate it on the largest public dataset OU-MVLP. As shown in Tab. 3, it can be seen that MetaGait outperforms other SOTA methods by considerable margins, which proves the generalizability of MetaGait.

#### 4.4 Ablation Study

This section presents ablation studies to validate the effectiveness of MTA and MTP, including the quantitative and qualitative analysis.

**Effectiveness of MTA and MTP.** The individual impacts of the MTA and MTP module are presented in Tab. 4. The baseline model refers to the feature extractor in [46] with traditional temporal aggregation (Max Pooling) and a separate FC layer. From the results, several conclusions are summarized as: 1) Using MTA or MTP individually can obtain **3.3%** and **2.5%** performance gain, respectively, which indicates the effectiveness of these modules. And MetaGait

Table 4: Ablation study on the effectiveness of each component of MetaGait, including Meta Triple Attention and Meta Temporal Aggregation.

Method	NM	BG	CL	Mean
Baseline	96.1	90.5	80.3	89.0
Baseline + MTA	97.5	94.2	85.4	92.3
Baseline + MTP	96.8	93.8	84.0	91.5
<b>MetaGait</b>	<b>98.1</b>	<b>95.2</b>	<b>86.9</b>	<b>93.4</b>

Table 5: Ablation study on the combination of receptive field of the local branch in omni-scale representation of Meta Triple Attention.

Local	NM	BG	CL	Mean
{1}	97.3	94.0	85.3	92.2
{1,3}	97.8	94.7	86.4	93.0
{1,3,5}	<b>98.1</b>	<b>95.2</b>	<b>86.9</b>	<b>93.4</b>
{1,3,5,7}	97.4	94.3	85.6	92.4
{1,3,5,7,9}	97.2	93.7	84.8	91.9

Table 6: Analysis of Meta Triple Attention, including the attention on three dimension, the calibration network, and the soft aggregation gate.

Attention			Calibration		Aggregation Gate	NM	BG	CL	Mean
Spatial	Channel	Temporal	Static	Meta					
✓			✓			96.8	93.8	84.0	91.5
	✓		✓			97.4	94.0	84.5	92.0
		✓	✓			97.0	94.1	84.1	91.7
			✓			96.8	94.0	84.7	91.8
✓	✓		✓			97.5	94.2	84.8	92.2
✓		✓	✓			97.6	94.3	85.0	92.3
	✓	✓	✓			97.1	94.1	84.9	92.0
✓	✓	✓	✓			97.7	94.5	85.2	92.5
✓	✓	✓	✓		✓	97.8	94.8	86.0	92.9
✓	✓	✓		✓		98.0	95.0	86.4	93.1
✓	✓	✓		✓	✓	<b>98.1</b>	<b>95.2</b>	<b>86.9</b>	<b>93.4</b>

improves the performance by **4.4%**; 2) Both MTA and MTP significantly improve the performance under the most challenging condition (*i.e.*, CL) by **5.1%** and **3.7%**. 3) The performance gain with MTP is mainly reflected in the BG/CL condition, where temporal aggregation would be more crucial [35].

**Receptive Field in Omni-scale Representation.** In MTA, we use the re-weighted combination of receptive fields in diverse scales to achieve omni-scale representation. To explore the effects of the different combinations, we use the convolutions with the kernel size of 1,3,5,7,9 in the local calibration network of MTA as shown in Tab. 5. It can be seen that the performance is improved with the increase of the receptive field scale until the combination of {1,3,5}. In contrast, a larger receptive field decreases the performance, which may lie in that larger and more diverse receptive fields could improve the ability of feature representation, but the over-parameterized convolution is hard to optimize.

**Analysis of MTA.** To evaluate the effectiveness of MTA, we analyze it from three aspects, *i.e.*, the attention design, the kind of the calibration network, and the soft aggregation gate. From the results shown in Tab. 6, we could conclude:

Table 7: The ablation study on Meta Temporal Aggregation.

Aggregation			Weight Network		NM	BG	CL	Mean
Max	Mean	GeM	Static	Meta				
✓			–	–	97.5	94.2	85.4	92.3
	✓		–	–	96.3	93.4	84.0	91.2
		✓	–	–	97.3	94.3	85.6	92.4
✓	✓		✓		97.3	94.2	85.7	92.4
	✓	✓	✓		97.6	94.5	85.7	92.6
✓		✓	✓		97.7	94.5	85.9	92.7
✓	✓	✓	✓		97.9	94.7	86.2	92.9
✓	✓	✓		✓	<b>98.1</b>	<b>95.2</b>	<b>86.9</b>	<b>93.4</b>



Fig. 6: The visualizaton of feature space using t-SNE [50].

1) MTA effectively improves the performance either using alone or in dimension combination; 2) The calibration network and the soft aggregation gate, which are parameterized by the meta-knowledge, clearly improve the rank-1 accuracy by **1.2%**. The above experimental results indicate that our MHN can effectively improve the model’s adaptiveness.

**Analysis of MTP.** The results in Tab. 7 shows the impacts of different temporal aggregation methods and weighting network. It can be seen that: 1) Different temporal aggregation methods used together provide performance gain by **0.6%**. 2) Weighting network with meta-knowledge could effectively integrate the merits of three aggregation methods than the static one, which indicates that MTP could achieve more comprehensive and discriminative representation.

**Visualization of Feature Space.** To validate the effectiveness of MetaGait intuitively, we randomly choose 10 IDs from CASIA-B to visualize their feature distribution. As shown in Fig. 6, we find that MetaGait improves the intra-class compactness and inter-class separability than baseline.

**Visualization of Attention Maps.** To qualitatively analyze MTA, we visualize the attention map shown in Fig. 7. For spatial dimension, MTA effectively learns the shape-aware attention map to guide the learning process adaptively. For temporal dimension, MTA can adaptively highlight important frames and suppress irrelevant frames to model the temporal representation. For channel dimension, it can be observed that MTA can learn a sample adaptive representation. Further, we can observe that different samples have low attention weights in certain channels, which may be caused by the channel redundancy in common.

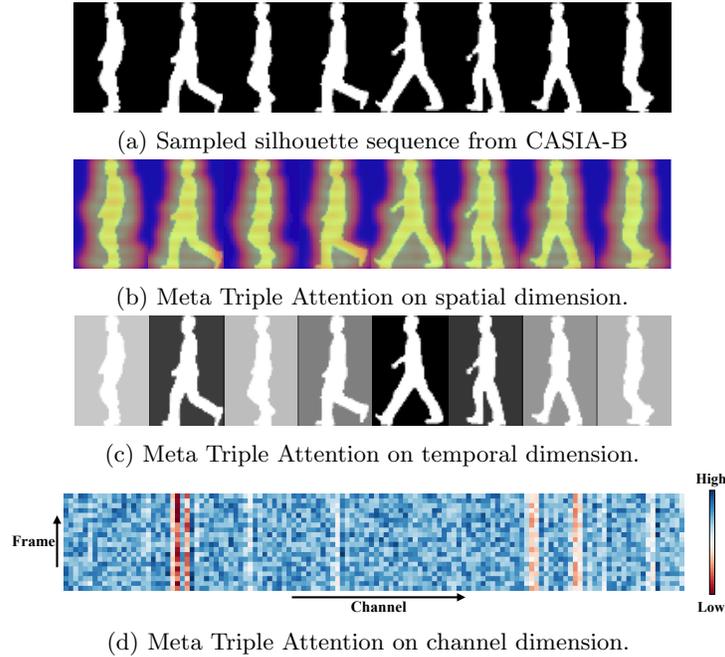


Fig. 7: The visualization of the attention maps of Meta Triple Attention. The transparency of the silhouette in (c) represents its attention value.

## 5 Conclusion

We propose a novel MetaGait framework to alleviate the conflicts between limited visual clues and various covariates with diverse scales. The key idea is to leverage meta-knowledge learned from Meta Hyper Network to improve the adaptiveness of attention mechanism. Specifically, Meta Triple Attention utilizes meta-knowledge to parameterize the calibration network and simultaneously conduct omni-scale attention on spatial/channel/temporal dimensions. Further, Meta Temporal Pooling excavates the relation between three complementary temporal aggregation methods and aggregates them in a sample adaptive manner. Finally, extensive experiments validate the effectiveness of MetaGait.

## Acknowledgements

This work is supported in part by the National Natural Science Foundation of China under Grant U20A20222, National Key Research and Development Program of China under Grant 2020AAA0107400, Zhejiang Provincial Natural Science Foundation of China under Grant LR19F020004, NSFC (62002320, U19B2043) and the Key R&D Program of Zhejiang Province, China (2021C01119).

## References

1. Antoniou, A., Edwards, H., Storkey, A.: How to train your maml. arXiv preprint arXiv:1810.09502 (2018)
2. Ariyanto, G., Nixon, M.S.: Model-based 3d gait biometrics. In: Int. Joint Conf. Bio. pp. 1–7 (2011)
3. Balazia, M., Plataniotis, K.N.: Human gait recognition from motion capture data in signature poses. IET Biom. pp. 129–137 (2017)
4. Bashir, K., Xiang, T., Gong, S.: Gait recognition using gait entropy image. In: IET Int. Conf. Imag. Crime Detect. Prevention. pp. 1–6 (2009)
5. Bengio, E., Bacon, P.L., Pineau, J., Precup, D.: Conditional computation in neural networks for faster models. arXiv preprint arXiv:1511.06297 (2015)
6. Bodor, R., Drenner, A., Fehr, D., Masoud, O., Papanikolopoulos, N.: View-independent human motion classification using image-based reconstruction. Int. Video Conf. pp. 1194–1206 (2009)
7. Bouchrika, I.: A survey of using biometrics for smart visual surveillance: Gait recognition. *Surveill. Action* pp. 3–23 (2018)
8. Boulgouris, N.V., Chi, Z.X.: Gait recognition based on human body components. In: IEEE Int. Conf. Image Process. pp. 353–356 (2007)
9. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estimation using part affinity fields. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 7291–7299 (2017)
10. Chao, H., He, Y., Zhang, J., Feng, J.: GaitSet: Regarding gait as a set for cross-view gait recognition. In: AAAI (2019)
11. Chattopadhyay, P., Sural, S., Mukherjee, J.: Frontal gait recognition from incomplete sequences using rgb-d camera. *IEEE Trans. Inf. Forensics Secur.* **9**(11), 1843–1856 (2014)
12. Chen, Y., Dai, X., Liu, M., Chen, D., Yuan, L., Liu, Z.: Dynamic convolution: Attention over convolution kernels. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 11030–11039 (2020)
13. Cheng, H.P., Zhang, T., Yang, Y., Yan, F., Li, S., Teague, H., Li, H., Chen, Y.: Swiftnet: Using graph propagation as meta-knowledge to search highly representative neural architectures. arXiv preprint arXiv:1906.08305 (2019)
14. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 764–773 (2017)
15. Denil, M., Shakibi, B., Dinh, L., Ranzato, M., De Freitas, N.: Predicting parameters in deep learning. *Adv. Neural Inform. Process. Syst.* (2013)
16. Devos, A., Chatel, S., Grossglauser, M.: Reproducing meta-learning with differentiable closed-form solvers. In: Int. Conf. Learn. Represent. (2019)
17. Dou, H., Zhang, P., Su, W., Yu, Y., Lin, Y., Li, X.: Gaitgci: Generative counterfactual intervention for gait recognition. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 5578–5588 (June 2023)
18. Dou, H., Zhang, P., Zhao, Y., Dong, L., Qin, Z., Li, X.: Gaitmpl: Gait recognition with memory-augmented progressive learning. *IEEE Trans. Image Process.* (2022)
19. Dou, H., Zhang, W., Zhang, P., Zhao, Y., Li, S., Qin, Z., Wu, F., Dong, L., Li, X.: Versatilegait: A large-scale synthetic gait dataset with fine-grained attributes and complicated scenarios. arXiv preprint arXiv:2101.01394 (2021)
20. Du, Y., Wang, W., Wang, L.: Hierarchical recurrent neural network for skeleton based action recognition. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 1110–1118 (2015)

21. Eigen, D., Ranzato, M., Sutskever, I.: Learning factored representations in a deep mixture of experts. arXiv preprint arXiv:1312.4314 (2013)
22. Fan, C., Peng, Y., Cao, C., Liu, X., Hou, S., Chi, J., Huang, Y., Li, Q., He, Z.: Gaitpart: Temporal part-based model for gait recognition. In: IEEE Conf. Comput. Vis. Pattern Recog. (2020)
23. Finn, C., Rajeswaran, A., Kakade, S., Levine, S.: Online meta-learning. In: Int. Conf. Mach. Learn. pp. 1920–1930 (2019)
24. Gao, H., Zhu, X., Lin, S., Dai, J.: Deformable kernels: Adapting effective receptive fields for object deformation. In: Int. Conf. Learn. Represent. (2019)
25. Goffredo, M., Bouchrika, I., Carter, J.N., Nixon, M.S.: Self-calibrating view-invariant gait biometrics. IEEE Trans. Cybern. **40**(4), 997–1008 (2009)
26. Guoying Zhao, Guoyi Liu, Hua Li, Pietikainen, M.: 3d gait recognition using multiple cameras. In: Int. Conf. Autom. Face Gesture Recog. pp. 529–534 (2006)
27. Han, J., Bhanu, B.: Individual recognition using gait energy image. IEEE Trans. Pattern Anal. Mach. Intell. **28**(2), 316–322 (2006)
28. Harley, A.W., Derpanis, K.G., Kokkinos, I.: Segmentation-aware convolutional networks using local attention masks. In: Int. Conf. Comput. Vis. pp. 5038–5047 (2017)
29. He, B., Yang, X., Wu, Z., Chen, H., Lim, S.N., Shrivastava, A.: Gta: Global temporal attention for video action understanding. arXiv preprint arXiv:2012.08510 (2020)
30. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 770–778 (2016)
31. Hermans, A., Beyer, L., Leibe, B.: In defense of the triplet loss for person re-identification. arXiv preprint arXiv:1703.07737 (2017)
32. Hospedales, T., Antoniou, A., Micaelli, P., Storkey, A.: Meta-learning in neural networks: A survey. arXiv preprint arXiv:2004.05439 (2020)
33. Hou, S., Cao, C., Liu, X., Huang, Y.: Gait lateral network: Learning discriminative and compact representations for gait recognition. In: Eur. Conf. Comput. Vis. pp. 382–398 (2020)
34. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 7132–7141 (2018)
35. Huang, X., Zhu, D., Wang, H., Wang, X., Yang, B., He, B., Liu, W., Feng, B.: Context-sensitive temporal feature learning for gait recognition. In: Int. Conf. Comput. Vis. pp. 12909–12918 (October 2021)
36. Huang, Y., Zhang, J., Zhao, H., Zhang, L.: Attention-based network for cross-view gait recognition. In: Adv. Neural Inform. Process. Syst. pp. 489–498 (2018)
37. Huang, Y., Wang, Y., Tai, Y., Liu, X., Shen, P., Li, S., Li, J., Huang, F.: Curricularface: Adaptive curriculum learning loss for deep face recognition. In: IEEE Conf. Comput. Vis. Pattern Recog. (2020)
38. Huang, Z., Xue, D., Shen, X., Tian, X., Li, H., Huang, J., Hua, X.S.: 3d local convolutional neural networks for gait recognition. In: Int. Conf. Comput. Vis. pp. 14920–14929 (October 2021)
39. Jacobs, R.A., Jordan, M.I., Nowlan, S.J., Hinton, G.E.: Adaptive mixtures of local experts. Neural computation **3**(1), 79–87 (1991)
40. Kastaniotis, D., Theodorakopoulos, I., Fotopoulos, S.: Pose-based gait recognition with local gradient descriptors and hierarchically aggregated residuals. J. Electron. Imaging **25**(6), 063019 (2016)
41. Kusakunniran, W., Wu, Q., Zhang, J., Li, H., Wang, L.: Recognizing gaits across views through correlated motion co-clustering. IEEE Trans. Image Process. **23**(2), 696–709 (2014)

42. Kusakunniran, W., Wu, Q., Zhang, J., Ma, Y., Li, H.: A new view-invariant feature for cross-view gait recognition. *IEEE Trans. Inf. Forensics Secur.* **8**(10), 1642–1653 (2013)
43. Li, S., Liu, W., Ma, H.: Attentive spatial-temporal summary networks for feature learning in irregular gait recognition. *IEEE Trans. Multimedia* **21**(9), 2361–2375 (2019)
44. Liao, R., Yu, S., An, W., Huang, Y.: A model-based gait recognition method with body pose and human prior knowledge. *Pattern Recog.* **98**, 107069 (2020)
45. Lin, B., Zhang, S., Bao, F.: Gait recognition with multiple-temporal-scale 3d convolutional neural network. In: *ACM Int. Conf. Multimedia*. pp. 3054–3062 (2020)
46. Lin, B., Zhang, S., Yu, X.: Gait recognition via effective global-local feature representation and local temporal aggregation. In: *Int. Conf. Comput. Vis.* pp. 14648–14656 (October 2021)
47. Lin, J., Rao, Y., Lu, J., Zhou, J.: Runtime neural pruning. *Adv. Neural Inform. Process. Syst.* (2017)
48. Lin, P., Sun, P., Cheng, G., Xie, S., Li, X., Shi, J.: Graph-guided architecture search for real-time semantic segmentation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* (June 2020)
49. Ma, J., Zhao, Z., Yi, X., Chen, J., Hong, L., Chi, E.H.: Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In: *SIGKDD*. pp. 1930–1939 (2018)
50. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *J. Mach. Learn. Res.* **9**(11) (2008)
51. Macoveciuc, I., Rando, C.J., Borrión, H.: Forensic gait analysis and recognition: standards of evidence admissibility. *J. Forensic Sci.* **64**(5), 1294–1303 (2019)
52. Makihara, Y., Sagawa, R., Mukaigawa, Y., Echigo, T., Yagi, Y.: Gait recognition using a view transformation model in the frequency domain. In: *Eur. Conf. Comput. Vis.* p. 151–163 (2006)
53. Radenović, F., Tolias, G., Chum, O.: Fine-tuning cnn image retrieval with no human annotation. *IEEE Trans. Pattern Anal. Mach. Intell.* **41**(7), 1655–1668 (2018)
54. Samangoei, S., Nixon, M.S.: Performing content-based retrieval of humans using gait biometrics. *Multimed. Tools Appl.* pp. 195–212 (2010)
55. Sepas-Moghaddam, A., Etemad, A.: Deep gait recognition: A survey. *arXiv preprint arXiv:2102.09546* (2021)
56. Shan, S., Li, Y., Oliva, J.B.: Meta-neighborhoods. *Adv. Neural Inform. Process. Syst.* **33**, 5047–5057 (2020)
57. Su, H., Jampani, V., Sun, D., Gallo, O., Learned-Miller, E., Kautz, J.: Pixel-adaptive convolutional neural networks. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 11166–11175 (2019)
58. Su, W., Miao, P., Dou, H., Wang, G., Qiao, L., Li, Z., Li, X.: Language adaptive weight generation for multi-task visual grounding. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 10857–10866 (June 2023)
59. Sun, K., Xiao, B., Liu, D., Wang, J.: Deep high-resolution representation learning for human pose estimation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* (2019)
60. Takemura, N., Makihara, Y., Muramatsu, D., Echigo, T., Yagi, Y.: Multi-view large population gait dataset and its performance evaluation for cross-view gait recognition. *IPSJ Trans. Comput. Vis. Appl.* **10**(1), 4 (2018)
61. Veit, A., Belongie, S.: Convolutional networks with adaptive inference graphs. In: *Eur. Conf. Comput. Vis.* pp. 3–18 (2018)

62. Wang, D., Zhang, S.: Unsupervised person re-identification via multi-label classification. In: IEEE Conf. Comput. Vis. Pattern Recog. (2020)
63. Wang, L., Ning, H., Tan, T., Hu, W.: Fusion of static and dynamic body biometrics for gait recognition. IEEE TCSVT **14**(2), 149–158 (2004)
64. Wang, X., Yu, F., Dou, Z.Y., Darrell, T., Gonzalez, J.E.: Skipnet: Learning dynamic routing in convolutional networks. In: Eur. Conf. Comput. Vis. pp. 409–424 (2018)
65. Woo, S., Park, J., Lee, J.Y., Kweon, I.S.: Cbam: Convolutional block attention module. In: Eur. Conf. Comput. Vis. pp. 3–19 (2018)
66. Yang, B., Bender, G., Le, Q.V., Ngiam, J.: Condconv: Conditionally parameterized convolutions for efficient inference. Adv. Neural Inform. Process. Syst. (2019)
67. Yu, S., Tan, D., Tan, T.: A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition. In: Int. Conf. Pattern Recog. pp. 441–444 (2006)
68. Zhang, F., Wah, B.W.: Supplementary meta-learning: Towards a dynamic model for deep neural networks. In: Int. Conf. Comput. Vis. pp. 4344–4353 (2017)
69. Zhang, P., Dou, H., Yu, Y., Li, X.: Adaptive cross-domain learning for generalizable person re-identification. In: Eur. Conf. Comput. Vis. pp. 215–232. Springer (2022)
70. Zhang, Y., Huang, Y., Yu, S., Wang, L.: Cross-view Gait Recognition by Discriminative Feature Learning. IEEE Trans. Image Process. (2019)