

# GaitRef: Gait Recognition with Refined Sequential Skeletons

Haidong Zhu\* Wanrong Zheng\* Zhaoheng Zheng Ram Nevatia  
University of Southern California

{haidongz|wanrongz|zhaoheng.zheng|nevatia@usc.edu}

## Abstract

Identifying humans with their walking sequences, known as gait recognition, is a useful biometric understanding task as it can be observed from a long distance and does not require cooperation from the subject. Two common modalities used for representing the walking sequence of a person are silhouettes and joint skeletons. Silhouette sequences, which record the boundary of the walking person in each frame, may suffer from the variant appearances from carried-on objects and clothes of the person. Framewise joint detections are noisy and introduce some jitters that are not consistent with sequential detections. In this paper, we combine the silhouettes and skeletons and refine the framewise joint predictions for gait recognition. With temporal information from the silhouette sequences, we show that the refined skeletons can improve gait recognition performance without extra annotations. We compare our methods on four public datasets, CASIA-B, OUMVLP, Gait3D and GREW, and show state-of-the-art performance.

## 1. Introduction

Gait recognition [6, 22, 30, 36] aims to find the uniqueness of the walking and posture sequence of a person, which has the advantage of being able to be acquired from long distance and without the subject’s cooperation. To recognize the gait sequence of a person, researchers have developed silhouettes-based methods, such as GaitSet [3], GaitPart [4] and GaitGL [17], and skeleton-based methods, such as GaitGraph [27]. However, both input modalities have some deficiencies. For binarized silhouettes, variations in clothes and carried-on objects, as shown in Figure 1 (a), introduce external ambiguity, while jitters in joint detection, as Figure 1 (b), decrease the skeleton accuracy.

In this paper, we introduce combination of silhouette sequences with skeletons and gaining the benefits of both modalities via refining the framewise skeletons with silhouette sequences. Since jitter in the detected skeletons are of a

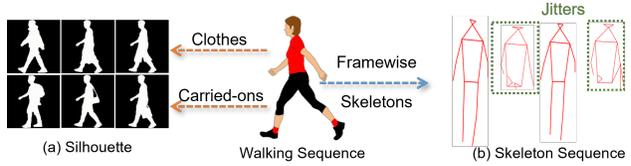


Figure 1. Visualization of the (a) silhouette and (b) skeleton sequence used for gait recognition. Silhouettes show different contours with different clothes and carried-on objects, while the skeletons suffer from jittery detection results in the video.

few frames isolated from the whole sequence, it is not temporally consistent with their neighbor frames [39]. Naive temporal smoothing, however, will introduce more confusion for gait recognition since the generated skeletons create new poses not consistent with the current sequence. Meanwhile, silhouettes for neighbor frames are of better temporal consistency due to the small changes in neighbor image conditions. We improve the quality of input skeletons by using silhouettes to fix the jitters while preserving necessary identity information for more precise gait recognition.

To combine and refine the silhouettes and skeletons, we introduce two methods, *GaitMix* and *GaitRef*. *GaitMix* takes skeletons and silhouettes as inputs and includes the encoded embeddings of both modalities for end-to-end gait recognition. *GaitMix* has two encoders, a silhouette feature encoder and a skeleton feature encoder, to project two modalities to their corresponding embedding spaces, followed by an MLP layer for fusing the concatenated feature to the identity embedding space. In our experiment, *GaitMix* works as the baseline since the existing state-of-the-art methods usually take only silhouette [3, 4, 17, 7] or skeleton [16, 27, 26] as the input of the network.

Based on the features encoded from *GaitMix*, *GaitRef* further refines the positions of the joints in the input skeleton sequence. We combine the encoded silhouette features with the encoded framewise skeleton features and the original joint positions to predict the relative changes for each point in the skeletons. Since the gait pattern should be consistent for the same person, features from the silhouettes and skeletons describing the same walking sequence

\* Equal contribution

should also be consistent, making the refinement of the skeletons with encoded silhouette features possible. In addition, the sequence-level silhouette feature helps the frame-level skeletons for each frame understand its corresponding poses without losing identity information since the temporal feature for the person is consistent and is shared across all the frames in the same walking sequence.

With the predicted change of the points, we add them back to the original skeleton sequence and use the skeleton encoder to extract the skeleton feature. We then concatenate it with the silhouette feature to predict the identity of the sequence with the refined skeletons. We assess our method on four public datasets, CASIA-B [37], OUMVLP [25], Gait3D [40] and GREW [44]. We show that the refined skeletons with silhouettes outperform other state-of-the-art gait recognition methods using skeletons and silhouettes as input, including *BaseMix*, our baseline method.

For the refinement of the input modalities, GaitEdge [15] introduced using RGB images to refine the silhouettes with the corresponding RGB images in the dataset. Due to privacy concerns, most public datasets [25, 40, 44] do not provide RGB images. We only require silhouettes and skeletons that are provided by the public datasets and achieve similar or even better results. We discuss more differences between the two methods in Sec. 4.1.

In summary, our contributions are 1) we introduce *GaitMix* and *GaitRef*, which combine the skeletons and silhouettes as end-to-end training for the gait recognition network, 2) we apply *GaitRef* for refining the skeletons with encoded silhouette features for refining skeletons without losing identity information in the sequence, and 3) we assess our model on four public datasets, CASIA-B, OUMVLP, Gait3D and GREW, and show state-of-the-art performance compared with other methods for gait recognition.

## 2. Related Work

**Gait Recognition** aims to find the corresponding identity of the person from the walking pattern. Considering the privacy issues in RGB images, gaits are usually recorded as two representations, silhouettes [37, 25] and skeletons [1]. Silhouettes record the boundary map of the human segmentation. To limit the impact of appearance variants on human shapes, researchers focus on part-based and body-shape reconstruction methods for gait recognition. GaitSet [3] and GLN [7] introduce set pooling and extract set features in the sequence. GaitPart [4] and GaitGL [17] split the image into different small patches and use local features to limit the impact of the appearance variants. In addition to directly mining identity information from silhouettes, ModelGait [13], Gait3D [40] and Gait-HBS [43] focus on 3-D shape reconstruction to assist the identification from sequences.

In addition to the mining identity from silhouette sequences, some researchers [16, 27] focus on using skeletons

instead of silhouettes for gait recognition. Compared with the body contours of the silhouettes, skeletons only include the joints and can remove the impact of body shapes as well as the appearance of the person. GaitGraph [27] uses the HRNet [28] for joint detections and uses the generated pose sequence for recognition. PoseGait [16] splits the gait sequence into pose, limb, angle, and motion, followed by analyzing the movements for each skeleton for these four features independently before combining them together for gait recognition. For the combination of silhouettes and skeletons, Wang *et al.* [29] directly concatenates the two features, which still suffers from erroneous joint detections.

**Pose Estimation and Refinement** focus on extracting the human body poses and refinement. With the development of transformers, pose estimation is also transforming from CNN-based networks [31, 2, 42] to transformer backbone networks [11, 14, 34, 32]. Pose estimation has experienced rapid development from CNNs [31] to vision transformer networks. Early works treat the transformer as a better decoder [11, 14, 34, 38]. Although the frame-level pose estimation accuracy is becoming more and more accurate, directly applying these methods to tasks with solid temporal relations, such as gait recognition, may introduce extra uncertainty with inaccurate joint predictions. For the sequence with strong temporal patterns, HuMoR [20] corrects the joint prediction of the person with the previous pose, and SmoothNet [39] filters the jitters in the whole sequence with analysis for the first and second deviation of the position for each point. These methods can fix some slight jitters in the long sequence but still suffer when the poses for a long sequence are inaccurate. For the task of gait recognition with temporal repeated patterns, even with inaccurate predictions for the long sequence, the model should still fix the joints with the consistent moving pattern of the same person, which these existing methods cannot achieve.

## 3. Methods

Given silhouettes  $S$  and joints  $J$  for the person  $p$ , the task of gait recognition is to match the identity with the people in a pool  $P = \{p_n\}_{n=1,2,\dots}$ , where  $n$  is the candidate identity. We encode  $S$  and  $J$  to their corresponding embeddings and find the nearest sample in  $P$  in the embedding space. In this section, we discuss the details of our proposed baseline *GaitMix*, which combines these two modalities, and the proposed method *GaitRef* to refine the input skeleton for gait recognition. We show both architectures in Figure 2.

In the remaining of this section, we first introduce *GaitMix* and *GaitRef* in Sec. 3.1 and 3.2, followed by the objectives for training in Sec. 3.3.

### 3.1. GaitMix: Multimodal Gait Recognition

*GaitMix* combines the skeletons and silhouettes as an end-to-end network for gait recognition. To extract the in-

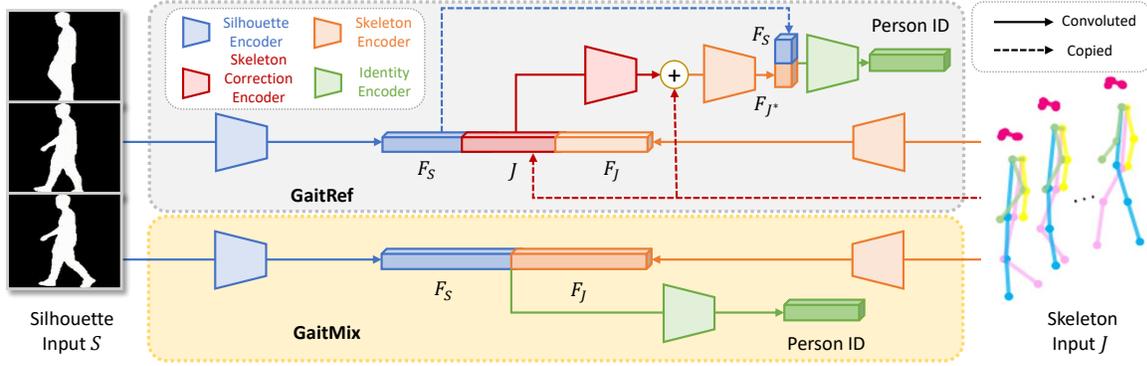


Figure 2. Our proposed architecture for GaitRef and GaitMix. Trapezoids are trainable modules, and modules of the same color in the same model share the weight. Dashed lines are the operation of feature copying.  $S$  and  $J$  are the input silhouettes and skeletons.  $F_S$  represents silhouette features, while  $F_J$  and  $F_{J^*}$  represent skeleton features from input and refined skeletons, respectively.

formation from both modalities, we apply two encoders: a silhouette feature encoder for encoding the silhouette  $S$ , and a skeleton feature encoder for projecting the input raw skeletons  $J$  into the embedding space.

**Silhouette Feature Encoder.** To extract the identity features from input sequential silhouette sequences, we use a silhouette feature encoder to convert the input silhouette sequence  $S$  to the corresponding output identity feature  $F_S$ . We have three steps for the silhouette feature encoder: convolution feature extraction, temporal pooling, and horizontal pooling. With the binary silhouette input sequence  $S = \{s_i\}_{i=1,\dots,N}$ , where  $i$  is the temporal stamp and  $N$  is the overall frame number, we apply a convolution network to extract the framewise feature  $f_i$  at frame  $i$ .  $f_i$  is an  $M$ -by- $N$ -by- $C$  matrix, where  $M$  and  $N$  are the height and width of the convoluted output features, and  $C$  is the channel number from the output of the last convolution layer.

With the framewise feature  $f_i$ , we use a max pooling layer for the temporal fusion and combine the feature into a single  $M$ -by- $N$ -by- $C$  output as temporal pooling. Since  $f_i$  still includes the spatial features for each segment, we follow [3, 5] and apply horizontal pyramid pooling with scale  $S$  as 5. The output of the feature is a  $2^{S-1}$ -by- $C$  feature vector after horizontal pooling. The architecture of each component can be found in the implementation details.

**Skeleton Feature Encoder.** In addition to the silhouette encoder, we have a skeleton feature encoder run in parallel and project the input skeleton sequence  $J$  to their corresponding human identification features  $F_J$ . For an input skeleton sequence  $J = \{j_i\}_{i=1,\dots,N}$ , where each input consists of  $K$  nodes and is shown as a  $K$ -by-2 matrix representing the 2-D skeletons for each frame, we follow [33] to apply spatial-temporal graph convolution network for the graphical feature extraction. By converting the input  $N$ -by- $K$ -by-2 to  $N$ -by- $K$ -by- $C$ , we average pool on the temporal and node dimensions and generate the final  $C$ -length vector as  $F_J$  representing the feature of the sequential skeleton.

**Fusion.** With the  $2^{S-1}$ -by- $C$  silhouette feature  $F_S$  and 1-by- $C$  skeleton feature  $F_J$  encoded from two different encoders, we concatenate the two features along with their first dimension and combine it to a  $(2^{S-1} + 1)$ -by- $C$  vector representing the body feature. We apply a shared MLP as an identity encoder for converting each  $C$ -length feature into the identity feature for identification.

### 3.2. GaitRef: Refining Skeletons with Silhouettes

Instead of directly combining skeleton and silhouette for gait recognition, *GaitRef* further uses the encoded feature from silhouette to improve the skeletons from the silhouette branch with temporal consistency. Since the errors in the skeleton generation are framewise jitters, temporal consistency can fix such jitters in the skeletons. In contrast, the refined skeletons can better help silhouettes ignore the appearance variants for gait recognition. Based on the architecture of *GaitMix*, *GaitRef* includes two external modules, a skeleton feature encoder which is shared with *GaitMix* pipeline and a skeleton correction network.

**Skeleton Correction Network.** With information from the silhouette feature, we use three different features as the network’s input to correct the skeleton and compute the corresponding adjustment for each point:  $2^{S-1}$ -by- $C$  silhouette features  $F_S$ ,  $N$ -by- $K$ -by- $C$  skeleton feature before pooling, and the original  $N$ -by- $K$ -by-2 joint matrix  $J$ .  $F_S$  provides the sequential information to correct the joint features  $F_J$ .  $F_J$  provides the framewise and feature for each node to correct the corresponding position of the joint in the frame.  $J$  provides the input order of the points to ensure the input and output order of the points are the same.

We show the architecture of the skeleton correction network in Figure 3. With these three inputs, we first flatten the silhouette feature into a  $2^{S-1} \times C$  vector. We then repeat it  $N$ -by- $K$  times and concatenate it with the other two features to form a  $N$ -by- $K$ -by- $(2^{S-1} \times C + C + 2)$  feature matrix. To decode the new position  $J'$  for each node

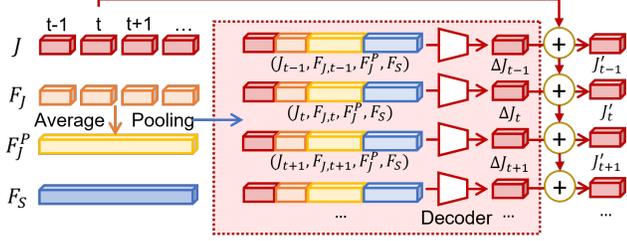


Figure 3. Architecture of the skeleton correction network.  $F_J^P$  is the skeleton features after average pooling. We concatenate the joint position  $J$  with its feature  $F_J$  along with the global feature after pooling  $F_J^P$  and the silhouette feature  $F_S$  before sending it into the decoder for calculating the position difference  $\Delta J$  for each frame. Decoders at different timestamps share weights.

in the sequence, we decode the  $\Delta J$  for all the points with a reversed spatial-temporal graph convolution network to decode the  $N$ -by- $K$ -by-2 adjustment for each node in  $J$ , and we have  $J'$  for refine the individual points in  $J$  following

$$J' = J + \Delta J = J + \text{SkeletonDecoder}(J, F_S, F_J) \quad (1)$$

The use of addition instead of directly predicting the corresponding location of the refined joints can give a relatively easier task for refinement and can preserve most of the original locations [41], since the original position of the joint has most of the sequential information correct and complete. By adding  $\Delta J$  on  $J$ , we get the final refined nodes as output and process it for further encoding.

**Skeleton Feature Encoder.** After we get the refined skeleton  $J'$  with  $\Delta J$ , we apply the same skeleton feature encoder used for *GaitMix* and apply it on the refined skeleton sequence  $J'$  for predicting the 1-by- $C$  skeleton feature  $F_{J'}$ . The two skeleton feature encoders share the parameters to ensure the two embedding spaces are the same between  $F_J$  and  $F_{J'}$ . Using the same skeleton feature encoder can also extend the data available for the encoder training to train a stabler graph convolution model for the feature extraction of the skeleton sequential input.

With the predicted  $F_{J'}$ , we concatenate it with  $2^{S-1}$ -by- $C$  silhouette feature  $F_S$  to form a  $(2^{S-1} + 1)$ -by- $C$  vector for representing the human body shape for *GaitRef*.

### 3.3. Objectives and Inference

We have two losses for both *GaitMix* and *GaitRef*. We use a triplet loss  $L_{triplet}$  for distinguishing the same identities in the same batch and a classification loss  $L_{cls}$  for the identities in training set with an MLP layer for projecting the identity feature to the number of candidates. For the combination of the two losses, we follow

$$L = \lambda_1 L_{triplet} + \lambda_2 L_{cls} \quad (2)$$

and empirically set  $\lambda_1$  as 1. For  $\lambda_2$  we follow [17, 40] to set it as different values for different datasets. We include

further discussion and the choice of parameters in the implementation details section in Sec. 4.1.

## 4. Experiments and Results

### 4.1. Experimental Details

**Datasets.** In our experiment, we assess our method on four public gait recognition datasets, CASIA-B [37], OUMVLP [6, 1], Gait3D [40] and GREW [44].

*CASIA-B* [37] has 124 subjects with 10 different walking variants for gait recognition. Among the 10 variants, 6 variants are for normal walking (NM), 2 variants are for the person carrying different bags (BG), and the remaining 2 variants are for different clothes (CL). Each subject has 110 videos captured with 10 variants from 11 different camera viewpoints distributed between  $0^\circ$  and  $180^\circ$ . We follow [3, 4, 7, 17] and use the videos of the first 74 identities for training and the remaining 50 for inference. During inference, we use the first four variances in normal walking conditions (NM) to build the gallery set as the library to query test sequences. The sequences of the remaining 2 NM variants, along with BG and CL sequences, are used as probe examples for finding the identity in the gallery.

*OUMVLP* [25, 1] is a large-scale dataset with 10,307 different identities. Each subject in this dataset has 2 different variants for normal walking (NM) conditions from 14 camera viewpoints, making 28 gait sequences. The angles of camera viewpoints are evenly distributed in two bins,  $0^\circ$  to  $90^\circ$  and  $180^\circ$  and  $270^\circ$ . Every two neighbor viewpoints have a 15-degree gap. We follow [3, 4, 7, 17] to use the identities with odd indexes between the 1-*st* and 10,305-*th* examples and build a training set with 5,153 identities. For the remaining 5,154 identities, we use the first sequence as the gallery set and the second as probes during inference.

*Gait3D* [40] is a medium dataset compared with CASIA-B and OUMVLP for gait recognition in the wild. It includes 4,000 identities among 25,309 video sequences captured via 39 cameras. Since sequences are captured in the wild, camera positions, carried-on objects, and clothes vary from sequence to sequence. Similar to GREW [44], *Gait3D* also provides both skeletons and silhouette sequences for each frame in the dataset. We follow [40] to use 3,000 identities for training and the remaining 1,000 during inference. For these 1,000 test cases, we build a probe set with 1,000 sequences for querying, as the probe set, and use the rest 5,369 sequences as the gallery set.

*GREW* [44] is a large in-the-wild gait recognition dataset with 128,671 sequences capturing 26,345 identities from 882 cameras. Each frame in the video has both silhouettes and poses provided. We follow [44] for using 20,000 identities for training and 6,000 identities as our test set. Each subject in the test set has 4 sequences, where we use two for the gallery and the other two as probes.

**Implementation Details.** For the implementation details section, we will discuss the details for the data preparation, model, and hyperparameter selection in experiments.

*Data preparation.* For all four datasets, we follow OpenGait<sup>1</sup> for preparing the silhouettes for each dataset and set the size of each frame as  $64 \times 44$ . Different from silhouettes, skeletons provided for different datasets are not exactly the same. Thus we process the skeletons for each dataset independently. For CASIA-B [37] dataset, we follow GaitGraph [27] and use a pretrained HR-Net [23] and generate the skeleton as MS COCO [18] format with 17 joints. The number of frames used for skeletons of CASIA-B is set to 60, and we use the 60 frames in the center of the whole sequence as our skeleton input.

For OUMVLP [25] dataset, we follow [1] for applying the skeletons along with the silhouette sequences, and we have skeleton sequences with 18 nodes per frame as OpenPose [2] format. Considering that the sequence length in OUMVLP is shorter than CASIA-B, we set the fixed frame number to 25 for each sequence. For videos shorter than 25, we repeat the frames until we have 25 frames.

For Gait3D [40] and GREW [44], since skeletons are collected in the wild, we normalize each skeleton by setting their height to 2 and move their center to the origin point  $(0, 0)$ . This can ensure that the position of the skeletons is aligned chiefly and will not change significantly.

*Network details.* In our network, we have two different encoders. For our silhouette feature encoder, we follow GaitGL [17] to build the encoder for CASIA-B, OUMVLP, and GREW. For Gait3D, we follow SMPLGait [40] and use its 2-D variant baseline, which we denote as OpenGait, to encode silhouette features. For the silhouette feature encoder in *GaitMix*, we follow ST-GCN [33] for encoding the skeletons into the same embedding dimension  $N_{out}$  as the silhouette feature encoder. The dimension of the hidden layers of ST-GCN is set to  $[64, 64, 128, 128, n_{out}]$ . In addition to the *GaitMix*, the decoder of the *GaitRef* uses the reversed shape of the ST-GCN, with  $[128, 64, 64, 3]$  as the hidden dimensions. For the encoder and decoder network, we have compared ST-GCN along with other choices, such as MS-G3D [19] in the ablation study.

*Model training.* In our model, we follow [17, 40] for choosing the hyperparameters. For CASIA-B, OUMVLP, and GREW, we use an Adam optimizer [10] with  $1e - 4$  as the learning rate for 80,000, 210,000, and 250,000 iterations, respectively. We decay the learning rate once at 70,000 iterations for CASIA-B and twice for OUMVLP and GREW, at iterations 150,000 and 200,000 as  $\frac{1}{10}$  of its original value. For the Gait3D dataset, we use the Adam optimizer for 180,000 iterations and set the initial learning rate as  $1e - 3$ , and the learning rate is decayed to  $\frac{1}{10}$  three times at iteration 30,000, 90,000 and 150,000. For CASIA-

B, OUMVLP and GREW, we follow [17] for using 1 for  $\lambda_1$  and  $\lambda_2$ , while we set  $\lambda_2$  as 0.1 for Gait3D following [40].

*Metrics and evaluations.* During inference, for each example in the probe set, we use  $L_2$  similarity to find the nearest example in the gallery set. For CASIA-B and OUMVLP, we evaluate the top-1 accuracy for the prediction. For GREW, we evaluate top-1, 5, 10 and 20 accuracies. For Gait3D, we assess top-1 and top-5 accuracies along with mAP and mINP following [35] for assessing since all the correct matches should have low-rank values when pairing the probe example with correct identities in the gallery.

**Baseline Methods.** For baseline methods, we compare with state-of-the-art gait recognition methods, including CNN-LB [30], GaitNet [22], GaitSet [3], GaitPart [4], GLN [7], GaitGL [17], ModelGait [13] and CSTL [8]. We also compare with PoseGait [16] and GaitGraph[27], which use skeleton sequences as the input. GaitEdge [15] generates silhouettes with RGB images<sup>2</sup> and MvModelGait [12] requires RGB images and camera positions, which are not provided by most of the datasets. Thus we do not make direct comparisons with them in our experiments.

## 4.2. Results and Analysis

To compare with other methods, we present both numerical results on gait recognition tasks for public datasets as well as the visualized generated skeletons from the refined branch, followed by the ablation studies.

**Numerical Results.** We show our numerical performance on the four datasets we used in Table 1, 2, 3 and 4. For CASIA-B and OUMVLP, identical-view cases are excluded. We have the following observations:

(i) *Comparison with other SOTA methods.* For all four different datasets we evaluate, we outperform the existing state-of-the-art methods with *GaitRef*. In Table 1, on CASIA-B, we achieve the best performance on all splits. Specifically, on NM, BG and CL, we reduce the error rates from 2.1%, 5.6% and 15.8% to 1.9%, 4.1%, 12.0%, which are relatively 6.7%, 26.8% and 25.3% reduction of the error rates. Even if we compare with GaitEdge [15] and MvModelGait [12], which use RGB images and viewpoint angles that do not usually exist in the public dataset, *GaitRef* still has a 1.6% and 7.3% improvement on CL, the hardest split.

For the other three datasets, on OUMVLP in Table 2, we ties with CSTL [8] for the top-1 accuracy, while we outperform it along with other methods for all the metrics on Gait3D [40] and GREW [44] in Table 3 and 4, which we show 2.7% and 1.6% improvements on Rank-1 accuracies respectively and consistent improvements on other metrics. This shows the solidness of correcting the skeleton using the knowledge in the silhouette sequence as *GaitRef*.

<sup>1</sup><https://github.com/ShiqiYu/OpenGait>

<sup>2</sup><https://github.com/ShiqiYu/OpenGait/tree/master/datasets/CASIA-B>

Probe	Method	Camera Positions											Mean
		0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°	
NM #5-6	PoseGait [16]	55.3	69.6	73.9	75.0	68.0	68.2	71.1	72.9	76.1	70.4	55.4	68.7
	CNN-LB [30]	83.3	92.3	96.7	94.6	91.7	89.7	92.2	94.0	96.3	92.3	79.0	91.1
	GaitNet [22]	91.2	92.0	90.5	95.6	86.9	92.6	93.5	96.0	90.9	88.8	89.0	91.6
	GaitGraph [27]	85.3	88.5	91.0	92.5	87.2	86.5	88.4	89.2	87.9	85.9	81.9	87.7
	GaitSet [3]	91.1	98.0	99.6	97.8	95.4	93.8	95.7	97.5	98.1	97.0	88.2	95.6
	GaitPart [4]	94.0	98.7	99.3	98.8	94.8	92.6	96.4	98.3	99.0	97.4	91.2	96.4
	GLN [7]	93.8	98.5	99.2	98.0	95.2	92.9	95.4	98.5	99.0	99.2	91.9	96.5
	GaitGL [17]	95.3	97.9	99.0	97.8	96.1	95.3	97.2	98.9	99.4	98.8	94.5	97.3
	CSTL [8]	97.2	99.0	99.2	98.1	96.2	95.5	97.7	98.7	99.2	98.9	96.5	97.8
	ModelGait [13]	96.9	97.1	98.5	98.4	97.7	98.2	97.6	97.6	98.0	98.4	98.6	<b>97.9</b>
	GaitMix	96.6	98.6	99.2	98.0	97.1	96.2	97.5	98.9	99.3	99.0	94.7	97.7
	GaitRef	97.2	98.7	99.1	98.0	97.3	97.0	98.0	99.4	99.4	98.9	96.4	<b>98.1</b>
	MvModelGait [12]	97.5	97.6	98.6	98.8	97.7	98.9	98.9	97.3	97.6	97.8	97.9	98.1
GaitEdge* [15]	97.2	99.1	99.2	98.3	97.3	95.5	97.1	99.4	99.3	98.5	96.4	97.9	
BG #1-2	PoseGait [16]	35.3	47.2	52.4	46.9	45.5	43.9	46.1	48.1	49.4	43.6	31.1	44.5
	CNN-LB [30]	64.2	80.6	82.7	76.9	64.8	63.1	68.0	76.9	82.2	75.4	61.3	72.4
	GaitNet [22]	83.0	87.8	88.3	93.3	82.6	74.8	89.5	91.0	86.1	81.2	85.6	85.7
	GaitGraph [27]	75.8	76.7	75.9	76.1	71.4	73.9	78.0	74.7	75.4	75.4	69.2	74.8
	GaitSet [3]	87.0	93.8	94.6	92.9	88.2	83.0	86.6	92.6	95.7	92.9	83.4	90.1
	GaitPart [4]	89.5	94.5	95.3	93.5	88.5	83.9	89.0	93.6	96.0	94.1	85.3	91.2
	GLN [7]	92.2	95.6	96.7	94.3	91.8	87.8	91.4	95.1	96.3	95.7	87.2	93.1
	GaitGL [17]	93.0	95.7	97.0	95.9	93.3	90.0	91.9	96.8	97.5	96.9	90.7	<u>94.4</u>
	CSTL [8]	91.7	96.5	97.0	95.4	90.9	88.0	91.5	95.8	97.0	95.5	90.3	93.6
	ModelGait [13]	94.8	92.9	93.8	94.5	93.1	92.6	94.0	94.5	89.7	93.6	90.4	93.1
	GaitMix	94.4	96.7	96.8	96.1	94.3	90.4	93.5	97.4	98.0	97.2	92.2	95.2
	GaitRef	94.4	96.4	97.3	96.8	96.2	92.2	94.4	97.2	98.7	97.9	93.3	<b>95.9</b>
	MvModelGait [12]	93.9	92.5	92.9	94.1	93.4	93.4	95.0	94.7	92.9	93.1	92.1	93.4
GaitEdge* [15]	95.3	97.4	98.4	97.6	94.3	90.6	93.1	97.8	99.1	98.0	95.0	96.1	
CL #1-2	PoseGait [16]	24.3	29.7	41.3	38.8	38.2	38.5	41.6	44.9	42.2	33.4	22.5	36.0
	CNN-LB [30]	37.7	57.2	66.6	61.1	55.2	54.6	55.2	59.1	58.9	48.8	39.4	54.0
	GaitNet [22]	42.1	58.2	65.1	70.7	68.0	70.6	65.3	69.4	51.5	50.1	36.6	58.9
	GaitGraph [27]	69.6	66.1	68.8	67.2	64.5	62.0	69.5	65.6	65.7	66.1	64.3	66.3
	GaitSet [3]	71.0	82.6	84.0	80.0	71.7	69.1	72.1	76.7	78.5	77.2	63.4	75.1
	GaitPart [4]	72.5	82.8	86.0	82.2	79.5	71.0	77.7	80.8	82.9	81.4	67.7	78.6
	GLN [7]	78.5	90.4	90.3	85.1	80.2	75.8	78.1	81.8	80.9	83.2	72.6	81.5
	GaitGL [17]	71.7	90.5	92.4	89.4	84.9	78.1	83.1	87.5	89.1	83.9	67.4	83.5
	CSTL [8]	78.1	89.4	91.6	86.6	82.1	79.9	81.8	86.3	88.7	86.6	75.3	<u>84.2</u>
	ModelGait [13]	78.2	81.0	82.1	82.8	80.3	76.9	75.5	77.4	72.3	73.5	74.2	77.6
	GaitMix	79.2	89.5	94.2	90.0	84.9	80.3	85.2	89.2	90.3	86.9	73.7	85.8
	GaitRef	81.4	93.3	94.3	91.6	87.8	83.9	88.5	91.7	91.6	89.1	75.0	<b>88.0</b>
	MvModelGait [12]	77.0	80.0	83.5	86.1	84.5	84.9	80.6	80.4	77.4	76.6	76.9	80.7
GaitEdge* [15]	84.3	92.8	94.3	92.2	84.6	83.0	83.0	87.5	87.4	85.9	75.0	86.4	

Table 1. Gait recognition results on CASIA-B dataset, excluding identical-view cases. GaitEdge\* requires RGB frames and uses the re-segmented CASIA-B\* silhouettes instead of CASIA-B, and MvModelGait requires the input camera viewpoints. We mark the best results among all the methods in bold and the best results in our baseline methods with underline.

(ii) *Comparison between GaitMix and GaitRef.* In addition to the comparison between the existing state-of-the-art methods, we also compare the performance between *GaitMix* and *GaitRef*, since these two methods use the skeleton modalities and are of the same setting for a fair comparison. We note that *GaitRef* outperforms *GaitMix* and shows pretty consistent improvements on all splits for all four datasets. While *GaitMix* introduces the skeleton information to discard the negative impact of the body contour, the inaccurate skeleton introduces some extra ambiguity, making the network unable to utilize the skeleton information maximally. With the refined skeletons, the model can capture more useful and accurate information for the identification task of the corresponding person in the video.

(iii) *Comparison with using 3-D body shapes.* Different from the other two datasets, Gait3D [40] provides the 3-D body shapes along with silhouette sequences, which are used by SMPLGait [40]. In Table 3 we provide the comparison for using 3-D body shapes as SMPLGait [40] and using skeletons as *GaitMix* and *GaitRef*. All these methods use OpenGait [40] as backbones. Compared with using silhouettes as the only input modality, the use of skeleton and body shape can both improve recognition accuracy. In SMPLGait, skeleton information is partially stored in the generated 3-D body shape for gait recognition, making *GaitMix* show similar performance as SMPLGait on all four metrics.

Compared with SMPLGait using 3-D body shape as the second modality, *GaitRef* with refined skeletons achieves a

Method	Camera Positions														Mean
	0°	15°	30°	45°	60°	75°	90°	180°	195°	210°	225°	240°	255°	270°	
GEINet [21]	23.2	38.1	48.0	51.8	47.5	48.1	43.8	27.3	37.9	46.8	49.9	45.9	45.7	41.0	42.5
GaitSet [3]	79.2	87.7	89.9	90.1	87.9	88.6	87.7	81.7	86.4	89.0	89.2	87.2	87.7	86.2	87.0
GaitPart [4]	82.8	89.2	90.9	91.0	89.7	89.9	89.3	85.1	87.7	90.0	90.1	89.0	89.0	88.1	88.7
GLN [7]	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 [17]	84.2	89.8	91.3	91.7	90.8	91.0	90.4	88.1	88.2	90.5	90.5	89.5	89.7	88.8	89.6
MvModelGait [12]	87.7	89.7	91.1	90.1	89.8	90.3	90.3	88.1	89.4	89.4	90.0	90.8	90.0	89.7	89.7
CSTL [8]	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	<b>90.2</b>
GaitMix	85.4	90.3	91.2	91.5	91.2	90.9	90.5	88.9	88.7	90.3	90.5	89.8	89.6	88.9	89.9
GaitRef	85.7	90.5	91.6	91.9	91.3	91.3	90.9	89.3	89.0	90.8	90.8	90.1	90.1	89.5	<b>90.2</b>

Table 2. Gait recognition results on OUMVLP dataset, excluding identical-view cases.

Methods	Rank@1	Rank@5	mAP	mINP
GaitSet [3]	36.70	58.30	30.01	17.30
GaitPart [4]	28.20	47.60	21.58	12.36
GLN [7]	31.40	52.90	24.74	13.58
GaitGL [17]	29.70	48.50	22.29	13.26
OpenGait [40]	42.90	63.90	35.19	20.83
CSTL [8]	11.70	19.20	5.59	2.59
SMPLGait [40]	<u>46.30</u>	<u>64.50</u>	<u>37.16</u>	<u>22.23</u>
GaitMix	45.80	65.60	36.74	22.09
GaitRef	<b>49.00</b>	<b>69.30</b>	<b>40.69</b>	<b>25.26</b>

Table 3. Gait recognition results reported on the Gait3D dataset with  $64 \times 44$  as input sizes. For all four metrics, higher values of the same metric indicate better performance.

Methods	Rank-1	Rank-5	Rank-10	Rank-20
PoseGait [16]	0.2	1.1	2.2	4.3
GaitGraph [27]	1.3	3.5	5.1	7.5
GEINet [21]	6.8	13.4	17.0	21.0
TS-CNN [30]	13.6	24.6	30.2	37.0
GaitSet [3]	46.3	63.6	70.3	76.8
GaitPart [4]	44.0	60.7	67.4	73.5
CSTL [8]	50.6	65.9	71.9	76.9
GaitGL [17]	<u>51.4</u>	<u>67.5</u>	<u>72.8</u>	<u>77.3</u>
GaitMix	52.4	67.4	72.9	77.2
GaitRef	<b>53.0</b>	<b>67.9</b>	<b>73.0</b>	<b>77.5</b>

Table 4. Rank-1, 5, 10 and 20 accuracies on GREW dataset.

better recognition performance. Considering that the generation of the SMPL body shapes also requires skeletons [24], inaccurate pose estimation in the 3-D body shape generation also makes the model difficult to understand the noisy body shapes with erroneous poses in SMPLGait [40], while *GaitRef* does not suffer this with refined skeletons.

**Ablation studies.** For ablation studies, we present two results on 1) different ways for combining the skeleton and silhouette features, 2) different skeleton encoder and decoder networks and comparison with other skeleton refinement methods, and 3) inputs of skeleton correction network. We have both experiments conducted on the CASIA-B dataset [37] for all three different settings and present the Top-1 accuracy for the final gait recognition results. We

Methods	Combination	NM	CL	BG
GaitGL	N/A	97.3	94.4	83.5
GaitMix	Padding	96.4	93.7	83.2
GaitMix	Concat.	<b>97.7</b>	<b>95.2</b>	<b>85.8</b>
GaitRef	Padding	97.5	94.6	85.8
GaitRef	Concat.	<b>98.1</b>	<b>95.9</b>	<b>88.0</b>

Table 5. Ablation results for different silhouette and skeleton feature combination on CASIA-B dataset for three splits. ‘Padding’ indicates the skeleton feature is padded on each of the feature of different scales, while ‘concat.’ means we concatenate the feature along with the scale dimension and use it only once.

Methods	Encoder	Decoder	NM	CL	BG
GaitMix	ST-GCN	N/A	97.7	95.2	85.8
GaitMix	MS-G3D	N/A	98.0	95.5	86.4
GaitRef	ST-GCN	ST-GCN	<b>98.1</b>	<b>95.9</b>	88.0
GaitRef	ST-GCN	MS-G3D	<b>98.1</b>	95.7	<b>88.5</b>
GaitRef	MS-G3D	ST-GCN	<b>98.1</b>	<b>95.9</b>	88.3
GaitMix	Average Smoothing		97.6	95.0	85.6
GaitMix	Gaussian Smoothing		97.7	95.2	85.9
GaitMix	SmoothNet [39]		97.4	94.4	83.8

Table 6. Ablations for different encoder and decoder combinations for *GaitMix* and *GaitRef* and different skeleton smoothing methods on CASIA-B datasets. Results are reported in Top-1 accuracy.

show the results in Table 5, 6 and 7 respectively and have the observations as follows:

(i) **Feature Combination.** In addition to concatenating the features, we also repeat and pad the skeleton feature along with each segment of the silhouette features, which we label as ‘padding’. We show the results in Table 5. For comparison, we also add the performance of GaitGL [17] in the table, which only uses the silhouette feature for gait recognition and is our backbone baseline on CASIA-B.

We note that for *GaitMix*, padding the skeleton feature along with each size of the silhouette feature has worse performance compared with the GaitGL baseline even if we use an external modality, while *GaitRef* has some minor improvements compared with GaitGL, indicating the raw skeleton sequence is relatively noisy and introduce exter-

Split	w/o $F_J$	w/o $F_J^P$	w/o $F_S$	Full SCN
NM	97.7	97.9	97.6	98.1
BG	95.4	95.9	95.3	95.9
CL	87.0	87.9	85.9	88.0

Table 7. Ablation results of different input for the skeleton correction network on CASIA-B. SCN is skeleton correction network.

nal ambiguity for gait recognition if concatenated with all scales of features without refinement. Nevertheless, concatenating the skeleton sequence only once, along with the silhouette, has better performance for both *GaitMix* and *GaitRef*. When the skeleton inputs do not dominate the feature input, the silhouette features can provide useful information from the noisy data before and after refinement.

(ii) **Encoder-decoder models.** For the choice of encoder-decoder, we choose between two of the state-of-the-art skeleton action recognition models, ST-GCN [33] and MS-G3D [19]. We show the results in Table 6. We show that both MS-G3D and ST-GCN can improve the performance of *GaitMix* and *GaitRef*. In our experiment, MS-G3D requires much larger GPU memory and at least  $\times 2$  training time for each introduced MS-G3D module. Considering the similar performance and time consumption, we use ST-GCN for both encoder and decoder networks.

(iii) **Different skeleton refinement methods.** For skeleton refinement, we compare *GaitRef* with neighbor smoothing (average and gaussian window for neighbor three frames) and SmoothNet [39] (pretrained on H36m [9]) on the CASIA-B dataset with Top-1 accuracy. We show the results in Table 6, where *GaitRef* outperforms other methods. For all three variations, 3-frame Gaussian smoothing has some small improvement compared with *GaitMix* but still cannot compete with *GaitRef*. Compared with naive temporal smoothing, which creates poses not consistent with the person in the sequence, combining the silhouette features introduces the walking patterns that do not exist in the skeletons and help them refine themselves for gait recognition. Compared with the refined skeletons from the skeleton sequence alone, external knowledge from the encoded silhouette embeddings eliminates ambiguity. It introduces ID-specific information during training when the walking pattern cannot be correctly extracted from the skeleton alone.

(iv) **Input of skeleton correction network.** Given the presence of three distinct inputs in our skeleton correction network besides  $J$ , namely  $F_J$ ,  $F_J^P$ , and  $F_S$ , we investigate the contributions of each component and present the results in Table 7 using three splits of CASIA-B datasets. We observe that when  $F_J$  and  $F_S$  are excluded, there is a significant performance drop, with  $F_S$  being the major contributor to the final correction. The skeleton correction network leverages temporal consistency in the skeleton sequences for correction, while the additional silhouette information offers external support for enhanced understanding.  $F_P$  has

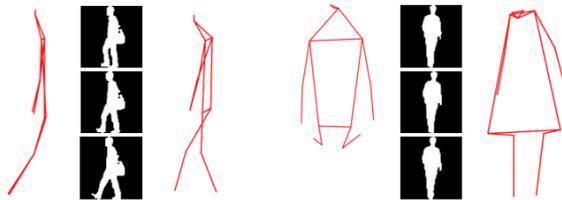


Figure 4. Visualization of successful and failure refined skeletons with *GaitRef*. For each example, from left to right, we have original skeletons, silhouette of the nearby timestamp and corrected skeletons from skeleton correction network.

limited utility, as it can be derived from  $F_J$ , while incorporating all three inputs leads to the best performance.

**Skeletons visualization.** We show two examples from the *GaitRef* compared to the original skeletons in Figure 4, accompanied by the three nearest silhouettes of similar time stamp. Through the *GaitRef* refinement, the resulting skeletons for gait recognition exhibit a considerable reduction in jitters and more accurately represent the individual’s walking patterns, especially the visibility of feet. Although the corrected skeletons may not be entirely precise on noisier datasets (such as CASIA-B), they improve over the initial jitters by using the same skeleton encoder to align the domain between input and refined skeletons and still positively contributes to the final recognition accuracy.

## 5. Conclusion

We introduce *GaitMix* and *GaitRef* for combining and refining the skeletons with silhouettes for gait recognition. *GaitMix* takes skeleton and gait sequences in an end-to-end network for projecting these two modalities into the same embedding space, while *GaitRef* further applies the temporal consistency in silhouettes for correcting the jitters in the skeletons. We show that combining the two modalities in *GaitMix* gives more accurate predictions, while the refined skeletons with silhouettes improve the quality of skeletons and generate more precise predictions. We assess our models on four public datasets, CASIA-B, OUMVLP, Gait3D, and GREW, and show state-of-the-art performance.

## Acknowledgement

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via [2022-21102100007]. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

## References

- [1] 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. *TBBIS*, 2(4):421–430, 2020. 2, 4, 5
- [2] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *TPAMI*, 43(1):172–186, 2019. 2, 5
- [3] H. Chao, Y. He, J. Zhang, and J. Feng. Gaitset: Regarding gait as a set for cross-view gait recognition. In *AAAI*, pages 8126–8133, 2019. 1, 2, 3, 4, 5, 6, 7
- [4] 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 *CVPR*, pages 14225–14233, 2020. 1, 2, 4, 5, 6, 7
- [5] 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 *AAAI*, volume 33, pages 8295–8302, 2019. 3
- [6] Y. He, J. Zhang, H. Shan, and L. Wang. Multi-task gans for view-specific feature learning in gait recognition. *TIFS*, 14(1):102–113, 2018. 1, 4
- [7] S. Hou, C. Cao, X. Liu, and Y. Huang. Gait lateral network: Learning discriminative and compact representations for gait recognition. In *ECCV*, pages 382–398, 2020. 1, 2, 4, 5, 6, 7
- [8] 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 *ICCV*, pages 12909–12918, 2021. 5, 6, 7
- [9] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI*, pages 1325–1339, 2013. 8
- [10] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [11] K. Li, S. Wang, X. Zhang, Y. Xu, W. Xu, and Z. Tu. Pose recognition with cascade transformers. In *CVPR*, pages 1944–1953, 2021. 2
- [12] X. Li, Y. Makihara, C. Xu, and Y. Yagi. End-to-end model-based gait recognition using synchronized multi-view pose constraint. In *ICCV*, pages 4106–4115, 2021. 5, 6, 7
- [13] X. Li, Y. Makihara, C. Xu, Y. Yagi, S. Yu, and M. Ren. End-to-end model-based gait recognition. In *ACCV*, 2020. 2, 5, 6
- [14] Y. Li, S. Zhang, Z. Wang, S. Yang, W. Yang, S.-T. Xia, and E. Zhou. Tokenpose: Learning keypoint tokens for human pose estimation. In *ICCV*, 2021. 2
- [15] J. Liang, C. Fan, S. Hou, C. Shen, Y. Huang, and S. Yu. Gaitedge: Beyond plain end-to-end gait recognition for better practicality. *arXiv preprint arXiv:2203.03972*, 2022. 2, 5, 6
- [16] R. Liao, S. Yu, W. An, and Y. Huang. A model-based gait recognition method with body pose and human prior knowledge. *PR*, 2020. 1, 2, 5, 6, 7
- [17] B. Lin, S. Zhang, and X. Yu. Gait recognition via effective global-local feature representation and local temporal aggregation. In *ICCV*, pages 14648–14656, 2021. 1, 2, 4, 5, 6, 7
- [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755, 2014. 5
- [19] Z. Liu, H. Zhang, Z. Chen, Z. Wang, and W. Ouyang. Disentangling and unifying graph convolutions for skeleton-based action recognition. In *CVPR*, pages 143–152, 2020. 5, 8
- [20] D. Rempe, T. Birdal, A. Hertzmann, J. Yang, S. Sridhar, and L. J. Guibas. Humor: 3d human motion model for robust pose estimation. In *ICCV*, pages 11488–11499, 2021. 2
- [21] K. Shiraga, Y. Makihara, D. Muramatsu, T. Echigo, and Y. Yagi. Geinet: View-invariant gait recognition using a convolutional neural network. In *ICB*, pages 1–8, 2016. 7
- [22] C. Song, Y. Huang, Y. Huang, N. Jia, and L. Wang. Gaitnet: An end-to-end network for gait based human identification. *PR*, 96:106988, 2019. 1, 5, 6
- [23] K. Sun, B. Xiao, D. Liu, and J. Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, pages 5693–5703, 2019. 5
- [24] Y. Sun, W. Liu, Q. Bao, Y. Fu, T. Mei, and M. J. Black. Putting people in their place: Monocular regression of 3d people in depth. In *CVPR*, pages 13243–13252, 2022. 7
- [25] 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. *TCVA*, 10(1):1–14, 2018. 2, 4, 5
- [26] T. Teepe, J. Gilg, F. Herzog, S. Hörmann, and G. Rigoll. Towards a deeper understanding of skeleton-based gait recognition. In *CVPRW*, pages 1569–1577, 2022. 1
- [27] T. Teepe, A. Khan, J. Gilg, F. Herzog, S. Hörmann, and G. Rigoll. Gaitgraph: Graph convolutional network for skeleton-based gait recognition. In *ICIP*, pages 2314–2318, 2021. 1, 2, 5, 6, 7
- [28] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, et al. Deep high-resolution representation learning for visual recognition. *TPAMI*, 43(10):3349–3364, 2020. 2
- [29] L. Wang, R. Han, J. Chen, and W. Feng. Combining the silhouette and skeleton data for gait recognition. *arXiv preprint arXiv:2202.10645*, 2022. 2
- [30] Z. Wu, Y. Huang, L. Wang, X. Wang, and T. Tan. A comprehensive study on cross-view gait based human identification with deep cnns. *TPAMI*, 39(2):209–226, 2016. 1, 5, 6, 7
- [31] B. Xiao, H. Wu, and Y. Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, 2018. 2
- [32] Y. Xu, J. Zhang, Q. Zhang, and D. Tao. Vitpose: Simple vision transformer baselines for human pose estimation. *arXiv preprint arXiv:2204.12484*, 2022. 2
- [33] S. Yan, Y. Xiong, and D. Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *AAAI*, 2018. 3, 5, 8
- [34] S. Yang, Z. Quan, M. Nie, and W. Yang. Transpose: Keypoint localization via transformer. In *ICCV*, 2021. 2
- [35] M. Ye, J. Shen, G. Lin, T. Xiang, L. Shao, and S. C. Hoi. Deep learning for person re-identification: A survey and outlook. *TPAMI*, 44(6):2872–2893, 2021. 5

- [36] S. Yu, H. Chen, Q. Wang, L. Shen, and Y. Huang. Invariant feature extraction for gait recognition using only one uniform model. *Neurocomputing*, 239:81–93, 2017. 1
- [37] S. Yu, D. Tan, and T. Tan. A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition. In *ICPR*, volume 4, pages 441–444, 2006. 2, 4, 5, 7
- [38] Y. Yuan, R. Fu, L. Huang, W. Lin, C. Zhang, X. Chen, and J. Wang. Hrformer: High-resolution transformer for dense prediction. In *NeurIPS*, 2021. 2
- [39] A. Zeng, L. Yang, X. Ju, J. Li, J. Wang, and Q. Xu. Smoothnet: A plug-and-play network for refining human poses in videos. *ECCV*, 2022. 1, 2, 7, 8
- [40] 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 *CVPR*, pages 20228–20237, 2022. 2, 4, 5, 6, 7
- [41] H. Zhu, Y. Yuan, Y. Zhu, X. Yang, and R. Nevatia. Open: Order-preserving pointcloud encoder decoder network for body shape refinement. In *ICPR*, 2022. 4
- [42] H. Zhu, Z. Zheng, and R. Nevatia. Temporal shift and attention modules for graphical skeleton action recognition. In *ICPR*, pages 3145–3151, 2022. 2
- [43] H. Zhu, Z. Zheng, and R. Nevatia. Gait recognition using 3-d human body shape inference. In *WACV*, pages 909–918, 2023. 2
- [44] Z. Zhu, X. Guo, T. Yang, J. Huang, J. Deng, G. Huang, D. Du, J. Lu, and J. Zhou. Gait recognition in the wild: A benchmark. In *ICCV*, pages 14789–14799, 2021. 2, 4, 5