

A Self-Supervised Gait Encoding Approach with Locality-Awareness for 3D Skeleton Based Person Re-Identification

Haocong Rao, Siqi Wang, Xiping Hu, Mingkui Tan, Yi Guo, Jun Cheng, Xinwang Liu, and Bin Hu

Abstract—Person re-identification (Re-ID) via gait features within 3D skeleton sequences is a newly-emerging topic with several advantages. Existing solutions either rely on hand-crafted descriptors or supervised gait representation learning. This paper proposes a *self-supervised* gait encoding approach that can leverage *unlabeled* skeleton data to learn gait representations for person Re-ID. Specifically, we first create self-supervision by learning to reconstruct unlabeled skeleton sequences reversely, which involves richer high-level semantics to obtain better gait representations. Other pretext tasks are also explored to further improve self-supervised learning. Second, inspired by the fact that motion's continuity endows adjacent skeletons in one skeleton sequence and temporally consecutive skeleton sequences with higher correlations (referred as *locality* in 3D skeleton data), we propose a locality-aware attention mechanism and a locality-aware contrastive learning scheme, which aim to preserve locality-awareness on intra-sequence level and inter-sequence level respectively during self-supervised learning. Last, with context vectors learned by our locality-aware attention mechanism and contrastive learning scheme, a novel feature named Contrastive Attention-based Gait Encodings (CAGEs) is designed to represent gait effectively. Empirical evaluations show that our approach significantly outperforms skeleton-based counterparts by 15-40% *Rank-1* accuracy, and it even achieves superior performance to numerous multi-modal methods with extra RGB or depth information. Our codes are available at <https://github.com/Kali-Hac/Locality-Awareness-SGE>.

Index Terms—Skeleton Based Person Re-Identification; Gait; Self-Supervised Deep Learning; Locality-Aware Attention; Contrastive Learning

1 INTRODUCTION

THE goal of person re-identification (Re-ID) is to re-identify the same person in a different scene or view. It plays a pivotal role in various applications like security authentication, human tracking, and role-based activity understanding [1]–[10]. To perform person Re-ID effectively, *gait* is one of the most useful human body clues, and it has aroused a growing interest in the research community since gait can be collected by unobtrusive methods without cooperative subjects [11]. Physiological and psychological studies [12], [13] reveal that human individuals behave differently when walking, and they are endowed with some relatively stable gait patterns (*e.g.*, stride length, angles of body joints). Such unique gait patterns are usually highly valuable for high-level tasks like gait recognition [14] and person Re-ID [15].

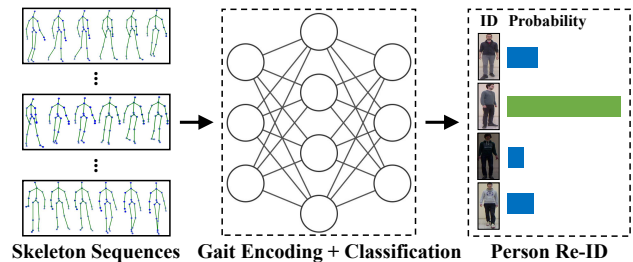


Fig. 1. Gait-based person Re-ID using 3D skeleton data.

To perform gait analysis, gait is typically described by two types of methods: (1) *Appearance*-based methods [16]–[23], which leverage human silhouettes from aligned image sequences to depict gait. However, an important flaw of this type of methods is its vulnerability to body shape changes and appearance changes. (2) *Model*-based methods [14], [15], [24]–[26], which model gait by human body structure and motion of body joints. Unlike appearance-based methods, model-based methods are invariant to factors like scale and view [27]. Therefore, model-based methods possess better robustness in practice. Among various models, *3D skeleton* model describes humans by the 3D coordinates of numerous key body joints, and it can often be used as a highly efficient representation of human body structure and motion [28]. 3D skeleton data are easily accessible with popular devices like Kinect [29], and they enjoy several prominent advantages when compared with frequently-seen RGB or depth data. For example, 3D skeleton data are much less likely to be interfered by illu-

- Haocong Rao and Siqi Wang contributed equally to this work (Corresponding authors: Xiping Hu; Bin Hu).
- H. Rao and J. Cheng are with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China. E-mail: {hc.rao, jun.cheng}@siat.ac.cn, haocongrao@gmail.com.
- S. Wang and X. Liu are with the National University of Defense Technology, Changsha 410073, China. Email: {wangsiqi10c, xinwangliu}@nudt.edu.cn.
- X. Hu is with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China, and also with Lanzhou University, Gansu 730000, China. E-mail: xp.hu@siat.ac.cn, huxp@lzu.edu.cn.
- M. Tan is with South China University of Technology, Guangzhou 510006, China. Email: mingkuitan@scut.edu.cn.
- Y. Guo is with the Second Clinical Medical College, Jinan University, Shenzhen 518055, China. E-mail: xuanyi_guo@163.com.
- B. Hu is with Beijing Institute of Technology, Beijing 100081, China, and also with Lanzhou University, Gansu 730000, China. Email: bh@bit.edu.cn, bh@lzu.edu.cn.

mination changes than RGB data, and they enjoy much smaller data size and less information redundancy than depth data [28]. Therefore, exploiting 3D skeleton data to perform gait analysis for downstream tasks like person Re-ID (illustrated in Fig. 1) is an attractive and promising solution with increasing popularity [14]. Nevertheless, the way to extract or learn discriminative gait features from 3D skeleton data remains to be an open problem.

For this purpose, most existing works like [15], [30] resort to hand-crafted skeleton feature descriptors. They typically focus on describing human bodies in terms of geometric, morphological, and anthropometric attributes, and then extracting corresponding features from 3D skeleton data. However, hand-crafted feature engineering is usually complicated and tedious. For instance, the method in [15] requires defining 80 skeleton descriptors from the views of anthropometric and gait attributes for person Re-ID. Besides, such methods also heavily rely on domain knowledge like human anatomy [31], thus lacking the ability to mine useful latent gait features beyond human comprehension. Motivated by the limitations of hand-crafted skeleton descriptors and the remarkable success achieved by recent deep neural networks (DNNs), few recent works like [32] resort to DNNs to learn gait representations automatically. However, the gait encoding process of such methods unexceptionally follows the classic *supervised* learning paradigm, which requires the discriminative information from labeled 3D skeleton data. As a consequence, they cannot utilize unlabeled skeleton data directly for automatic gait encoding.

This paper for the first time proposes a 3D skeleton based person Re-ID approach guided by *self-supervision* and *locality-awareness*, and it realizes highly effective gait encoding with unlabeled 3D skeleton sequence data. By first creating self-supervision signals for gait encoding, our approach not only makes it possible to learn gait representations from unlabeled skeleton data, but also prompts learning richer high-level semantics (*e.g.*, sequence order, body part motion) and more discriminative gait features. To be more specific, we propose to leverage the reverse reconstruction of skeleton sequences as a primary self-supervised learning objective. Meanwhile, we also explore other pretext tasks and utilize them to further enhance self-supervised learning. Second, we notice that 3D skeleton sequences are endowed with a property named *locality*: The continuity of human motion usually induces very small pose/skeleton changes in a local temporal interval [33]. As a result, for each skeleton in a skeleton sequence, its adjacent skeletons have higher correlations to itself, which is referred as *intra-sequence* locality. Similarly, we also define *inter-sequence* locality, which suggests that two temporally consecutive 3D skeleton sequences also enjoy higher relevance. To this end, we propose to incorporate *locality-awareness* to enable better 3D skeleton reconstruction and gait encoding during self-supervised learning. Accordingly, during the gait encoding process, we propose a novel locality-aware attention mechanism and locality-aware contrastive learning scheme to preserve locality on the intra-sequence and inter-sequence level respectively. Last, based on the proposed locality-aware attention mechanism and locality-aware contrastive learning, we devise a novel method to construct our gait representations, which are named Contrastive Attention-based Gait Encodings (CAGEs), from the learned model. Our empirical evaluations demonstrate that CAGEs, which can be learned without any skeleton label, can be directly applied to person Re-ID and achieve highly competitive performance.

A preliminary version of this work was reported in [34]. Compared with [34], this work not only systematically explores

the design of self-supervised learning objective for 3D skeleton sequences with more pretext tasks, but also extends the conception of locality from the intra-sequence level to the inter-sequence level by devising the locality-aware contrastive learning. To our knowledge, this is also the first attempt that explores the contrastive learning technique for learning discriminative gait features. In particular, we need to point out that “self-supervised learning” in this paper still refers to learning with our designed pretext tasks (*e.g.*, reverse reconstruction or prediction of the 3D skeleton sequences). Although contrastive learning is a popular technique to realize self-supervised learning in the literature, it is specifically designed to encourage sequence-level locality here and should be distinguished from the previous “self-supervised learning” term. On the foundation of those improvements, this work also improves earlier gait features proposed in [34] into the more effective gait features CAGEs. To validate those improvements, this work carries out more extensive experiments and detailed discussions on three public Re-ID datasets and a new multi-view Re-ID dataset [35]. Besides, we demonstrate that our approach is also effective with the skeleton data estimated from RGB videos [36]. To sum up, our contributions can be summarized as follows:

- We propose a new self-supervised learning paradigm for the gait encoding of 3D skeleton based person Re-ID. The proposed paradigm enables us to yield more effective gait representations from unlabeled 3D skeleton sequences by learning a reverse sequential skeleton reconstruction.
- We explore other possible forms of pretext tasks for the proposed self-supervised learning paradigm, and showcase their effectiveness in further strengthening gait encoding.
- We devise a locality-aware attention mechanism to exploit the intra-sequence locality within skeletons of one skeleton sequence, so as to facilitate better skeleton reconstruction and gait encoding during the self-supervised learning.
- We propose a locality-aware contrastive learning scheme to preserve the inter-sequence locality among temporally adjacent 3D skeleton sequences, which is able to encourage better gait encoding on the sequence level.
- We propose a new method to construct our gait representations (CAGEs) from the learned model. CAGEs are shown to be highly effective for person Re-ID.

The rest of paper is organized as follows: Sec. 2 introduces relevant works in the literature. Sec. 3 elucidates each module of the proposed approach. Sec. 4 presents the details of experiments, and extensively compares our approach with existing solutions. Sec. 5 provides ablation studies and comprehensive discussions on the proposed approach. Sec. 6 concludes this paper.

2 RELATED WORKS

2.1 Person Re-identification

Skeleton-based person Re-ID. As an emerging topic, most existing works extract hand-crafted features to depict certain geometric, morphological or anthropometric attributes of 3D skeleton data. Barbosa *et al.* [30] compute 7 Euclidean distances between the floor plane and joint or the joint pair to construct a distance matrix, which is learned by a quasi-exhaustive strategy to perform person Re-ID. Munaro *et al.* [37] further extend them to 13 skeleton descriptors (D^{13}) and use support vector machine (SVM) and k -nearest neighbor (KNN) for classification. In [38], 16 Euclidean distances between body joints (D^{16}) are fed to an

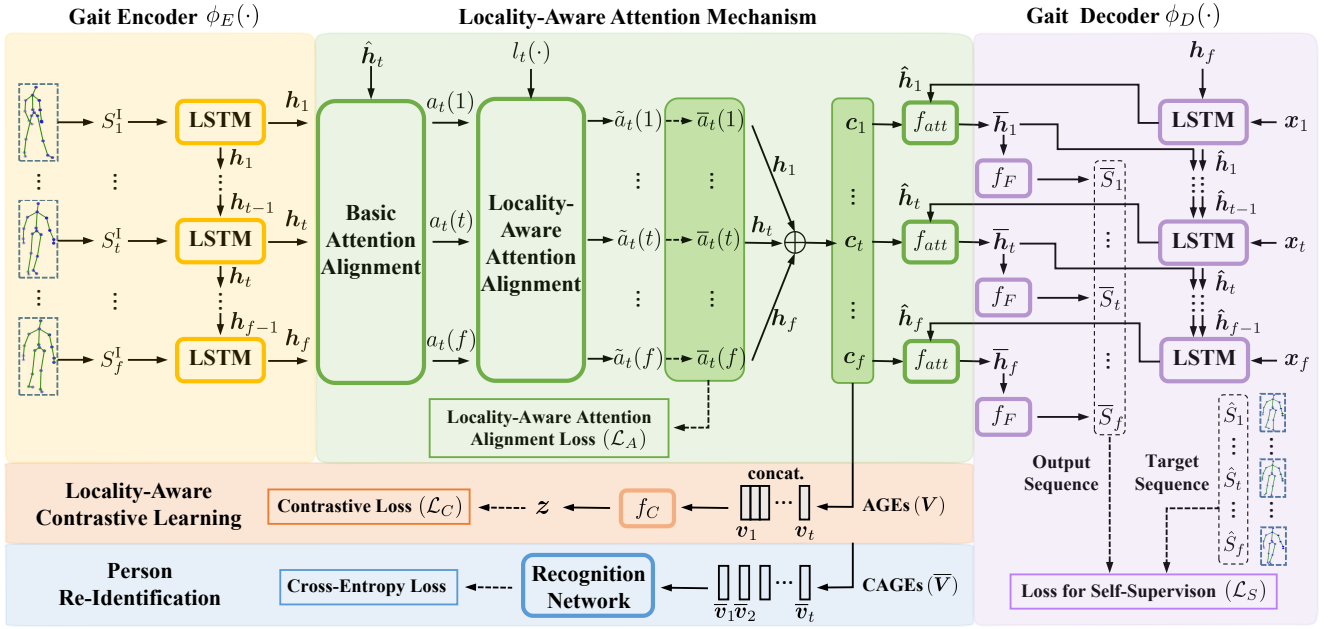


Fig. 2. Flow diagram of our model: (1) Gait Encoder (yellow) encodes each skeleton frame S_t^i into an encoded gait state h_t . (2) Locality-aware attention mechanism (green) first computes the basic attention alignment score $a_t(\cdot)$, so as to measure the content-based correlation between each encoded gait state and the decoded gait state \hat{h}_t from Gait Decoder (purple). Then, the locality mask $l_t(\cdot)$ provides an objective $\tilde{a}_t(\cdot) = a_t(\cdot) l_t(\cdot)$, which guides our model to learn locality-aware alignment scores $\bar{a}_t(\cdot)$ by the locality-aware attention alignment loss \mathcal{L}_A . Next, $h_1 \dots h_f$ are weighted by $\bar{a}_t(\cdot)$ to compute the context vector c_t . c_t and \hat{h}_t are fed into the concatenation layer $f_{att}(\cdot)$ to produce an attentional state vector \bar{h}_t . Finally, \bar{h}_t is fed into the full connected layer $f_F(\cdot)$ to output t^{th} skeleton \bar{S}_t and Gait Decoder for later decoding. (3) c_t is used to build Attention-based Gait Encodings (AGEs) v_t , which are concatenated and fed into $f_C(\cdot)$ to perform locality-aware contrastive learning (red). (4) The Contrastive Attention-based Gait Encodings (CAGEs) \bar{v}_t learned by contrastive learning are fed into a recognition network for person Re-ID (blue).

Adaboost classifier for Re-ID. Since existing solutions that use features from 3D skeletons alone usually perform unsatisfactorily, features from other modalities (*e.g.*, 3D point clouds [39], 3D face descriptor [38]) are often used to enhance the performance. Meanwhile, few recent works exploit supervised deep learning models to learn gait representations from skeleton data: In [40], a Time based Graph (TG) LSTM model is proposed for human recognition based on skeletal graphs that are transformed from binary images. [32] utilizes long short-term memory (LSTM) [41] to model temporal dynamics of body joints to perform person Re-ID; The latest work from Liao *et al.* [14] propose PoseGait, which feeds 81 hand-crafted pose features of 3D skeleton data into convolutional neural networks (CNN) for human recognition.

Our work differs from previous skeleton-based works in following aspects: (1) We propose a novel self-supervised approach to encode discriminative gait features from unlabeled 3D skeleton data. We do NOT need to extract hand-crafted features like [30], [37], [38] or use identity labels to supervise gait representation learning [14], [32], [40]. In this work, the reverse skeleton reconstruction is proposed as the major pretext task to capture high-level semantics like skeleton motion patterns in unlabeled skeleton data, which facilitates us to yield more effective gait representations. (2) The property of locality induced by motion's continuity is exploited for better gait encoding: We propose the locality-aware attention mechanism and locality-aware contrastive learning scheme to preserve intra-sequence and inter-sequence locality embedded in 3D skeleton sequences respectively. To our best knowledge, this is also the first attempt to leverage attention mechanism and contrastive learning to realize gait encoding.

Depth-based and multi-modal person Re-ID. Depth-based methods typically exploit human shapes or silhouettes from depth

images to extract gait features for person Re-ID. For example, Sivapalan *et al.* [21] extend the Gait Energy Image (GEI) [20] to 3D domain and propose Gait Energy Volume (GEV) algorithm based on depth images to perform gait-based human recognition. 3D point clouds from depth data are also pervasively used to estimate body shape and motion trajectories. Munaro *et al.* [37] propose point cloud matching (PCM) to compute the distances of multi-view point cloud sets, so as to discriminate different persons. Haque *et al.* [32] adopt 3D LSTM to model motion dynamics of 3D point clouds for person Re-ID. As to multi-modal methods, they usually combine skeleton-based features with extra RGB or depth information (*e.g.*, depth shape features based on point clouds [39], [42], [43]) to boost Re-ID performance. In [44], CNN-LSTM with reinforced temporal attention (RTA) is proposed for person Re-ID based on a split-rate RGB-depth transfer approach.

2.2 Contrastive Learning

Recent years witness a surging popularity of contrastive learning in the unsupervised learning field [45]–[48]. It aims to learn effective data representations by separating positive pairs from negative pairs with contrastive losses, which measure the similarity of sample pairs in a latent representation space, and they are often combined with various pretext tasks to enhance unsupervised learning. To name a few, Wu *et al.* [46] propose an instance-level discrimination method based on exemplar related task [49] and noise-contrastive estimation (NCE) [50]. Contrastive predictive coding (CPC) [47] adopts the context auto-encoding task with a probabilistic contrastive loss (InfoNCE) to learn representations from different modalities. Many previous works [48], [51], [52] adopt the memory bank [46] to store the representation vectors of samples in the dataset, while some recent advances [53]–[55]

explore the use of in-batch samples for negative sampling instead of a memory bank. The latest contrastive learning framework is SimCLR [55], [56], which is highly efficient for unsupervised visual representation learning and inspires our work for skeletons.

Our work differs from previous studies in the following aspects: (1) The proposed locality-aware contrastive learning scheme is proposed to incorporate the inter-sequence locality into the gait encoding process, during which the sequence-level representation of 3D skeleton sequence is viewed as an instance in contrastive learning. (2) The goal of locality-aware contrastive learning scheme is to maximize the agreement between adjacent sequences that enjoy higher correlations. Different from [55], [57] that use augmented samples of images as contrastive instances, we exploit consecutive and non-consecutive 3D skeleton sequences as positive and negative pairs respectively for contrastive learning.

3 THE PROPOSED APPROACH

Suppose that a skeleton sequence $\mathbf{S} = (S_1, \dots, S_f)$ contains f consecutive skeleton frames, where $S_i \in \mathbb{R}^{J \times 3}$ contains 3D coordinates of J body joints. The training set $\Phi = \{\mathbf{S}^{(i)}\}_{i=1}^N$ contains N skeleton sequences collected from different persons. Each skeleton sequence $\mathbf{S}^{(i)}$ corresponds to a label y_i , where $y_i \in \{1, \dots, C\}$ and C is the number of persons. Our goal is to learn discriminative gait features \mathbf{v} from \mathbf{S} without using any label. Then, the effectiveness of learned features \mathbf{v} can be validated by using them to perform person Re-ID: Learned features and labels are used to train a simple recognition network (note that learned features \mathbf{v} are frozen and NOT tuned by training at this stage). The overview of the proposed approach is given in Fig. 2, and we present details of each technical component below.

3.1 Self-Supervised Learning with 3D Skeletons

3.1.1 Reverse Reconstruction as Self-Supervision

To learn gait representations without labeled 3D skeleton sequences, we propose to introduce self-supervision by learning to reconstruct input 3D skeleton sequences in a *reverse* order, *i.e.*, by taking the input skeleton sequence $\mathbf{S}^I = (S_1^I, \dots, S_f^I) = (S_1, \dots, S_f) = \mathbf{S}$, we expect our model to output the sequence $\hat{\mathbf{S}} = (\hat{S}_1, \dots, \hat{S}_f) = (S_f, \dots, S_1)$, which gives $\hat{S}_t = S_{f-t+1}$. Compared with the naïve reconstruction that learns to reconstruct exact inputs ($\mathbf{S}^I \mapsto \mathbf{S}^I$), the proposed learning objective ($\mathbf{S}^I \mapsto \hat{\mathbf{S}}$) is combined with more high-level information (*e.g.*, skeleton order in the sequence) that are meaningful to human perception, which requires the model to capture richer high-level semantics to achieve this learning objective. In this way, our model is expected to learn more meaningful gait representations than frequently-used plain reconstruction. Formally, given an input 3D skeleton sequence, we use the encoder to encode each skeleton frame S_t^I ($t \in \{1, \dots, f\}$) and the previous step's latent state \mathbf{h}_{t-1} (if existed), which provides the temporal context information for the gait encoding process, into the current latent state \mathbf{h}_t :

$$\mathbf{h}_t = \begin{cases} \phi_E(S_1^I) & \text{if } t = 1 \\ \phi_E(\mathbf{h}_{t-1}, S_t^I) & \text{if } t > 1 \end{cases} \quad (1)$$

where $\mathbf{h}_t \in \mathbb{R}^K$, $\phi_E(\cdot)$ denotes our Gait Encoder (GE). GE is built with an LSTM, which aims to capture long-term temporal dynamics of skeleton sequences. $\mathbf{h}_1, \dots, \mathbf{h}_f$ denote the *encoded gait states* that contain preliminary gait encoding information. In the *training* phase, encoded gait states are decoded by a Gait

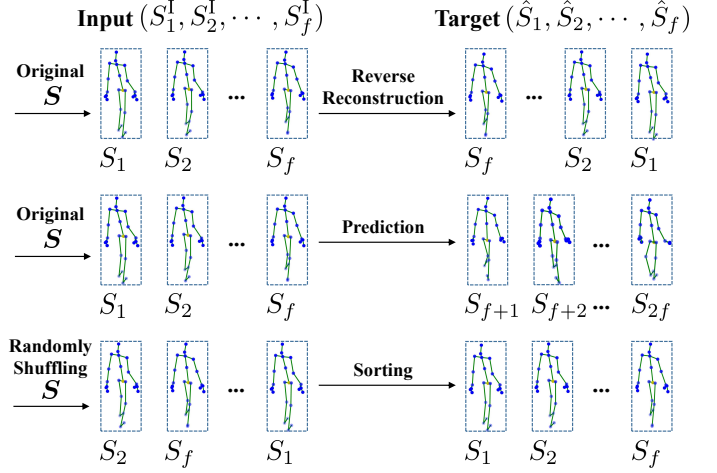


Fig. 3. Schematic diagrams of three pretext tasks: Reverse reconstruction (top), prediction (middle), sorting (bottom). The original sequence $\mathbf{S} = (S_1, \dots, S_f)$ is the input (S_1^I, \dots, S_f^I) for reverse reconstruction and prediction, and a random shuffle of \mathbf{S} is the input for sorting.

Decoder (GD) that aims to output the target sequence $\hat{\mathbf{S}}$, and the decoding process is performed below (see Fig. 2):

$$(\hat{\mathbf{h}}_t, \bar{S}_t) = \begin{cases} \phi_D(\mathbf{h}_f, \mathbf{x}_1) & \text{if } t = 1 \\ \phi_D(\hat{\mathbf{h}}_{t-1}, \mathbf{x}_{t-1}, \bar{\mathbf{h}}_{t-1}) & \text{if } t > 1 \end{cases} \quad (2)$$

where $\phi_D(\cdot)$ denotes the GD. GD consists of an LSTM and a fully connected (FC) layer that outputs those joint coordinates of a 3D skeleton. $\hat{\mathbf{h}}_t \in \mathbb{R}^K$ refers to the t^{th} *decoded gait state*, *i.e.*, the latent state output by GD's LSTM to generate the t^{th} skeleton \bar{S}_t . \mathbf{x}_t is the t^{th} auxiliary input for *training*. When the decoding is initialized ($t = 1$), we feed $\mathbf{x}_1 = \mathbf{Z} \in \mathbb{R}^J$, which is an all-0 skeleton placeholder, and the final encoded gait state \mathbf{h}_f into GD to decode the first skeleton. Afterwards, to generate t^{th} skeleton \bar{S}_t , $\phi_D(\cdot)$ takes three inputs from the $(t-1)^{\text{th}}$ decoding step: decoded gait state $\hat{\mathbf{h}}_{t-1}$, the auxiliary skeleton input $\mathbf{x}_{t-1} = \hat{S}_{t-1}$ (the ground-truth skeleton of the previous time step) that enables better convergence, and the *attentional state vector* $\bar{\mathbf{h}}_{t-1}$ that fuses encoding and decoding information based on the proposed attention mechanism, which will be elaborated in Sec. 3.2. In this way, we define the objective function \mathcal{L}_S for self-supervision, which minimizes the mean square errors (MSE) between a target sequence $\hat{\mathbf{S}}$ and an output sequence $\bar{\mathbf{S}}$:

$$\mathcal{L}_S = \sum_{i=1}^f \sum_{j=1}^J (\bar{S}_{ij} - \hat{S}_{ij})^2 \quad (3)$$

where $\bar{S}_{ij}, \hat{S}_{ij}$ represent the j^{th} joint position of the i^{th} output or target skeleton. In the *testing* phase, to test the reconstruction ability of our model, it should be noted that we use the output skeleton \bar{S}_{t-1} rather than the target skeleton \hat{S}_{t-1} as the auxiliary input to ϕ_D in the $t > 1$ case, *i.e.*, $\mathbf{x}_{t-1} = \bar{S}_{t-1}$. To facilitate training, our implementation actually optimizes Eq. 3 on each individual dimension of the skeleton's 3D coordinates: $\mathbf{S}^{[d]} \mapsto \hat{\mathbf{S}}^{[d]}$, where $d \in \{X, Y, Z\}$ corresponds to a certain dimension of the 3D data space, and $\mathbf{S}^{[d]}, \hat{\mathbf{S}}^{[d]} \in \mathbb{R}^{f \times J}$.

3.1.2 Other Pretext Tasks for Self-Supervision

Our self-supervised gait encoding approach can also be equipped with other pretext tasks, which exploit different inputs \mathbf{S}^I and learning targets $\hat{\mathbf{S}}$ to provide self-supervision for gait encoding of 3D skeleton sequences. To this end, we design two additional pretext tasks in this work: (1) Future skeleton frame prediction (“Prediction”). As shown in Fig. 3, the prediction task takes the original skeleton sequence as the input, namely $(S_1^I, \dots, S_f^I) = (S_1, \dots, S_f)$, and the learning goal is to predict the next f skeleton frames: $(\hat{S}_1, \dots, \hat{S}_f) = (S_{f+1}, \dots, S_{2f})$. The motivation of this task is that the model must capture key motion patterns in a skeleton sequence to predict unseen future skeletons, and learning to acquire such inference ability enables the model to mine more latent gait features. (2) Skeleton sequence sorting (“Sorting”). The sorting task attempts to sort a randomly shuffled 3D skeleton sequence \mathbf{S} back to the original sequence. Specifically, the input is $(S_1^I, \dots, S_f^I) = (S_{r_1}, \dots, S_{r_f})$ ($r_1, \dots, r_f \in \{1, \dots, f\}$ are shuffled indexes), and the target sequence is $(\hat{S}_1, \dots, \hat{S}_f) = (S_1, \dots, S_f)$. In this way, it enables the model to learn the inherent temporal coherence embedded in skeleton sequences during gait encoding. Besides, it is easy to know that reverse reconstruction is a special case of sorting. As a comparison, sorting is usually more difficult for the model, since a random shuffle often removes the sequence order information completely. By contrast, reverse reconstruction still retains the sequence order information at inputs, which makes it possible to utilize the locality property discussed in Sec. 3.2.

For those new pretext tasks of self-supervised learning, we can still leverage the same model structure while changing the inputs \mathbf{S}^I , targets $\hat{\mathbf{S}}$, and auxiliary input \mathbf{x} accordingly during *training*: For prediction and sorting, we alternate the auxiliary input from $\mathbf{x}_t = \hat{\mathbf{S}}_t$ (the ground-truth skeleton) to $\mathbf{x}_t = \bar{\mathbf{S}}_t$ (the predicted skeleton) in the $t > 1$ case, and in the testing phase we keep the auxiliary inputs unchanged. For prediction, we also test the half-prediction case (*i.e.*, target sequence is $\hat{\mathbf{S}} = (S_{\frac{f}{2}+1}, \dots, S_{\frac{3f}{2}})$), which is shown to yield better gait representations for person Re-ID than the full-prediction (see supplementary material). Our later experiments compare the gait features learned by different pretext tasks for person Re-ID, and the results demonstrate that the proposed reverse reconstruction achieves the best performance (see Sec. 4.3), which explains the center role of reverse reconstruction in the proposed self-supervised gait encoding approach. However, our experiments also show that gait features learned by other pretext tasks can be readily combined with features learned by reverse reconstruction (referred as the “Rev. Rec.” configuration in later experiments) to further improve Re-ID performance. First, we extensively evaluate different combinations of two pretext tasks: “Rev. Rec. + Pred.” (reverse reconstruction+prediction), “Rev. Rec. + Sort.” (reverse reconstruction+sorting), “Pred. + Sort.” (prediction+sorting) in Sec. 5.2.1. Then, we propose the “Rev. Rec. Plus” configuration that synthesizes gait features learned from all three pretext tasks for person Re-ID. Consequently, the exploration of more specific pretext tasks will be beneficial to our self-supervised skeleton sequence learning paradigm.

3.2 Locality-Aware Attention Mechanism

As learning gait features essentially requires capturing motion patterns from 3D skeleton sequences, it is instinctive to consider a natural property of motion–continuity. The continuity of motion ensures that those skeletons in a small temporal interval will

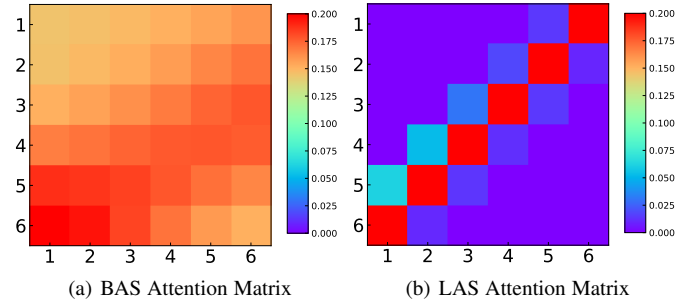


Fig. 4. Visualization of the BAS (left) and LAS (right) attention matrices that represent average attention alignment scores. Note that the abscissa and ordinate denote indices of input skeletons and output skeletons respectively. The LA alignment improves the learning of locality by assigning larger alignment scores near the clinodiagonal line.

NOT undergo drastic changes, thus resulting in higher correlations among adjacent skeletons in a local context of the skeleton sequence. This property is referred as *intra-sequence locality* here. Due to such intra-sequence locality, when reconstructing a certain skeleton in a sequence, we expect our model to pay more attention to its neighboring skeletons that are located in the same local temporal context. To this end, we propose a locality-aware attention mechanism, the details of which are presented below.

3.2.1 Basic Attention Alignment

We first introduce Basic Attention (BA) alignment [58] to measure the content (*i.e.*, latent state) based correlations between the input sequence and the output sequence. As shown in Fig. 2, at the t^{th} decoding step, we compute the *BA Alignment Scores (BAS)* between $\hat{\mathbf{h}}_t$ and the encoded gait state \mathbf{h}_j ($j \in \{1, \dots, f\}$):

$$a_t(j) = \text{align}(\hat{\mathbf{h}}_t, \mathbf{h}_j) = \frac{\exp(\hat{\mathbf{h}}_t^\top \mathbf{h}_j)}{\sum_{i=1}^f \exp(\hat{\mathbf{h}}_t^\top \mathbf{h}_i)} \quad (4)$$

BAS aims to focus on those more correlative skeletons in the encoding stage, and provides preliminary attention weights for skeleton decoding. However, BA alignment only considers the content based correlations and does not explicitly take the intra-sequence locality into consideration, which motivates us to design the locality mask and locality-aware attention alignment below.

3.2.2 Locality Mask

The motivation to design locality mask is to incorporate intra-sequence locality directly into the gait encoding process for better skeleton reconstruction. As the goal is to decode the t^{th} skeleton $\bar{\mathbf{S}}_t$ as $\hat{\mathbf{S}}_t$, we consider those skeletons in the local temporal context of S_{p_t} to be highly correlated to $\bar{\mathbf{S}}_t$, where $p_t = f - t + 1$ (note that we use the reverse reconstruction). To describe the local context centered at S_{p_t} , we define an attentional window $[p_t - D, p_t + D]$, where D is a selected integer to control the attentional range. Since the locality will favor temporal positions near p_t (*i.e.*, closer positions are more correlative), a direct solution is to place a *Gaussian* distribution centered around p_t as the locality mask:

$$l_t(j) = \exp\left(-\frac{(j - p_t)^2}{2\sigma^2}\right) \quad (5)$$

where we empirically set $\sigma = \frac{D}{2}$, j is a position within the window centered at p_t . We can weight BAS by this locality mask

to compute *Masked BA Alignment Scores (MBAS)* below, which directly forces alignment scores to obtain locality:

$$\tilde{a}_t(j) = l_t(j) \cdot a_t(j) \quad (6)$$

Besides, the locality mask is only valid for sequential reconstruction, so it cannot be directly combined with sorting or prediction. This is exactly an advantage of the reverse reconstruction task.

3.2.3 Locality-Aware Attention Alignment

Despite that the locality mask is straightforward to yield the intra-sequence locality, it is a very coarse solution that brutally constrains the alignment scores. Therefore, instead of using MBAS ($\tilde{a}_t(j)$) directly, we propose the *Locality-aware Attention (LA) alignment*. Specifically, an LA alignment loss term \mathcal{L}_A is used to encourage LA alignment to learn similar locality like $\tilde{a}_t(j)$:

$$\mathcal{L}_A = \sum_{t=1}^f \sum_{j=1}^f (a_t(j) - \tilde{a}_t(j))^2 \quad (7)$$

By adding the loss term \mathcal{L}_A , we can obtain *LA Alignment Scores (LAS)*. Note that in Eq. 4, the final learned $a_t(j)$ is BAS. For clarity, we use $\bar{a}_t(j)$ to represent LAS learned by Eq. (7). With the guidance of \mathcal{L}_A , our model learns to allocate more attention to the local temporal context by itself rather than using a hard locality mask. To utilize alignment scores to yield an attention-weighted encoded gait state at the t^{th} step, we can calculate the *context vector* \mathbf{c}_t by a sum of weighted encoded gait states:

$$\mathbf{c}_t = \sum_{j=1}^f \bar{a}_t(j) \mathbf{h}_j \quad (8)$$

Note that the context vector \mathbf{c}_t can also be computed with BAS or MBAS. \mathbf{c}_t provides a synthesized gait encoding that is more relevant to $\hat{\mathbf{h}}_t$, which facilitates the reconstruction of t^{th} skeleton. To combine both encoding and decoding information for reverse reconstruction, we use a concatenation layer $f_{att}(\cdot)$ that combines \mathbf{c}_t and $\hat{\mathbf{h}}_t$ into an attentional state vector $\bar{\mathbf{h}}_t$:

$$\bar{\mathbf{h}}_t = f_{att}(\mathbf{c}_t, \hat{\mathbf{h}}_t) = \tanh(\mathbf{W}_{att}[\mathbf{c}_t; \hat{\mathbf{h}}_t]) \quad (9)$$

where \mathbf{W}_{att} represents the learnable weight matrix in the layer. Finally, we generate the joint coordinates of t^{th} output skeleton by the FC layer $f_F(\cdot)$ of the GD:

$$\bar{\mathbf{s}}_t = f_F(\bar{\mathbf{h}}_t) = \mathbf{W}_F \bar{\mathbf{h}}_t \quad (10)$$

where \mathbf{W}_F is the weights to be learned in this FC layer.

3.2.4 Analysis on Different Attention Mechanisms

First, to provide a more intuitive impression of the proposed locality-aware attention mechanism, we visualize the BAS and LAS attention matrices, which are formed by the average alignment scores computed with the $f = 6$ skeleton sequence on dimension X (explained in Sec. 3.1.1) as an example. As shown by Fig. 4, LA alignment significantly improves the locality of learned alignment scores: It can be observed that relatively large alignment scores are densely distributed near the clinodiagonal line of the attention matrix (note that when reverse reconstruction is performed, the clinodiagonal line reflects each skeletons' correlations to themselves), which means temporally adjacent skeletons are assigned with larger attention than those comparatively remote skeletons. By contrast, despite that BA alignment also learns locality to a certain extent, BA's alignment weights exhibit a much

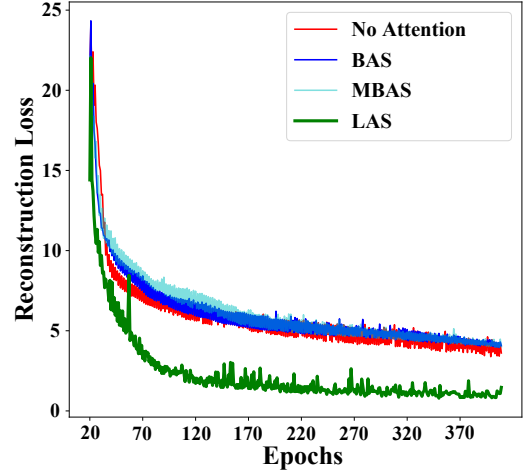


Fig. 5. Reconstruction loss curves when using no attention, BAS, MBAS or LAS for skeleton reconstruction (note that here we compare different attention mechanisms without using contrastive learning). Using LAS achieves better reverse reconstruction with smaller reconstruction loss.

more uniform distribution and many non-adjacent skeletons are also given large alignment scores. Similar trends are also observed when learning on dimension Y and Z . Such observations justify LAS's effectiveness to encourage intra-sequence locality.

Second, to illustrate how different attention mechanisms contribute to the self-supervised learning goal, *i.e.* the reverse reconstruction of 3D skeleton sequences, we visualize the corresponding reconstruction loss during training in four cases: No attention mechanism, BAS, MBAS and LAS. As shown by Fig. 5, it can be observed that training with LAS converges at a faster speed with a smaller reconstruction loss, which justifies our intuition that exploiting locality will facilitate the reverse reconstruction. Interestingly, we observe that using the locality mask directly in fact does not benefit the reduction of reconstruction loss, which also indicates that learning is a better way to accomplish intra-sequence locality than imposing a hard locality mask.

3.2.5 Attention-based Gait Encodings

Instead of simply fulfilling the pretext tasks, the ultimate goal here is to learn good gait features to conduct effective person Re-ID. Thus, we need to extract certain 3D skeleton sequence embeddings from the internal layers of neural networks to construct gait representations. Unlike traditional LSTM based methods that basically rely on the last hidden state to compress the temporal dynamics of a sequence [59], we recall that the dynamic context vector \mathbf{c}_t learned from the attention mechanism integrates the key encoded gait states of input skeletons and retains crucial spatio-temporal information to recover target skeleton sequences. Hence, we utilize them instead of the last hidden state to build the preliminary gait representations—Attention-based Gait Encodings (AGEs). Specifically, skeleton-level AGE (\mathbf{v}_t) is defined as follows:

$$\mathbf{v}_t = [\mathbf{c}_t^X; \mathbf{c}_t^Y; \mathbf{c}_t^Z] \quad (11)$$

where \mathbf{c}_t^d denotes the context vector computed on dimension $d \in \{X, Y, Z\}$ at the t^{th} step of decoding. As reported in our earlier work [34], AGEs can be directly utilized to perform person Re-ID. However, AGEs only incorporate the locality on the intra-sequence level, *i.e.*, among different skeletons in one 3D skeleton sequence. They do not consider the relationship between different

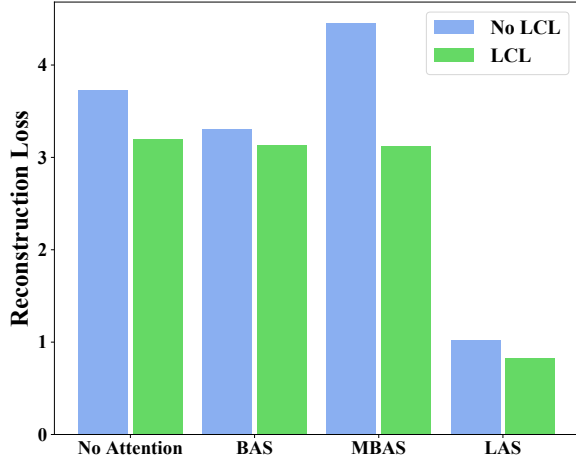


Fig. 6. Skeleton reconstruction loss when using no attention, BAS, MBAS or LAS for skeleton reconstruction. The comparison between applying LCL (“LCL”) and not applying LCL (“No LCL”) is reported.

3D skeleton sequences on the inter-sequence level, which can also be involved to improve the gait representation learning. This motivates us to propose the Locality-aware Contrastive Learning (LCL) scheme below, so as to further encourage the inter-sequence locality among different skeleton sequences. Besides, since LCL is still performed on each individual dimension and AGEs will be further tuned by LCL, we use a slightly abused notation by defining $\mathbf{v}_t = \mathbf{c}_t^d$ in the next section, where $d \in \{X, Y, Z\}$.

3.3 Locality-Aware Contrastive Learning Scheme

Similar to intra-sequence locality among skeletons in a sequence, we can also assume that consecutive skeleton sequences in a local temporal context are more likely to share similar gait representations than those non-consecutive ones. Such relationship among different 3D skeleton sequences (referred as *inter-sequence locality*) can be exploited to enhance self-supervised gait encoding. To this end, we propose a Locality-aware Contrastive Learning (LCL) scheme to impose inter-sequence locality on skeleton sequences.

3.3.1 Skeleton Sequence Contrastive Learning

To compute the correlations between skeleton sequences and learn the inter-sequence locality, we first construct sequence-level gait representations by concatenating skeleton-level AGEs as follows:

$$\mathbf{V}^{(i)} = [\mathbf{v}_1; \mathbf{v}_2; \dots; \mathbf{v}_t] \quad (12)$$

where $\mathbf{V}^{(i)}$ denotes attention-based gait encodings of the i^{th} skeleton sequence $\mathbf{S}^{(i)}$ in the training set Φ . Here we adopt the same setting in [55] to improve representation learning: We first use a multi-layer perceptron (MLP) with one hidden layer to map $\mathbf{V}^{(i)}$ to the contrasting space where the contrastive loss is applied:

$$\mathbf{z}_i = f_C(\mathbf{V}^{(i)}) = \mathbf{W}^2 \sigma(\mathbf{W}^1 \mathbf{V}^{(i)}) \quad (13)$$

where $\mathbf{z}_i \in \mathbb{R}^K$ is the representation of i^{th} skeleton sequence in the contrasting space. $f_C(\cdot)$ is the function that denotes the MLP layer for contrastive learning. σ is the non-linear activation function like ReLU. \mathbf{W}^1 and \mathbf{W}^2 are weights to be learned in the MLP layer. The LCL scheme contrasts the similarity between representations (\mathbf{z}_i) of different skeleton sequences. During the training stage of LCL scheme, each batch of skeleton sequences

Algorithm 1 Main algorithm of LCL scheme

Input: Batch size n , temperature τ , gait encoding model ϕ , MLP function f_C for contrastive learning.

for a batch of consecutive sequences $\{\mathbf{S}^{(k)}\}_{k=1}^n$ **do**

for all $k \in \{1, \dots, n-1\}$ **do**

$\mathbf{V}^{(k)} = \phi(\mathbf{S}^{(k)})$ # gait representation

$\mathbf{z}_k = f_C(\mathbf{V}^{(k)})$ # map to contrasting space

$\mathbf{V}^{(k+1)} = \phi(\mathbf{S}^{(k+1)})$ # $\mathbf{S}^{(k+1)}$ and $\mathbf{S}^{(k)}$ are adjacent

$\mathbf{z}_{k+n-1} = f_C(\mathbf{V}^{(k+1)})$

end for

for all $i \in \{1, \dots, 2n-2\}$ and $j \in \{1, \dots, 2n-2\}$ **do**

$\alpha_{i,j} = \frac{\mathbf{z}_i^\top \mathbf{z}_j}{\tau \|\mathbf{z}_i\| \|\mathbf{z}_j\|}$ # similarity between sequences

end for

define $\ell(i, j) = -\log \frac{\exp(\alpha_{i,j})}{\sum_{k=1}^{2n-2} \mathbb{1}_{[k \neq i]} \exp(\alpha_{i,k})}$

$\mathcal{L}_C = \frac{1}{2n-2} \sum_{k=1}^{n-1} [\ell(k, k+n-1) + \ell(k+n-1, k)]$

update ϕ and f_C to minimize \mathcal{L}_C

end for

return gait encoding model ϕ and f_C

$\{\mathbf{S}^{(k)}\}_{k=1}^n$ is drawn without random shuffle from the training subset Φ_i ($i \in \{1, \dots, C\}$) that corresponds to the i^{th} person, and $\Phi = \{\Phi_1, \dots, \Phi_C\}$. Next, we define two consecutive skeleton sequences as a positive pair, while two non-consecutive skeleton sequences in the batch are defined as a negative pair. Given a positive skeleton sequence pair $\mathbf{S}^{(i)}$ and $\mathbf{S}^{(j)}$ from $\{\mathbf{S}^{(k)}\}_{k=1}^n$, the LCL scheme aims to maximize the agreement between representations of $\mathbf{S}^{(i)}$ and $\mathbf{S}^{(j)}$. We define the loss function below for a positive sequence pair and summarize the entire LCL scheme in Algorithm 1:

$$\ell(i, j) = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j) / \tau)}{\sum_{k=1}^{2n-2} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)} \quad (14)$$

where $\text{sim}(\mathbf{z}_i, \mathbf{z}_j) = \mathbf{z}_i^\top \mathbf{z}_j / \|\mathbf{z}_i\| \|\mathbf{z}_j\|$ denotes the cosine similarity between two representation vectors \mathbf{z}_i and \mathbf{z}_j , τ denotes the temperature parameter, and $\mathbb{1}_{[k \neq i]} \in \{0, 1\}$ is an indicator function: $\mathbb{1}_{[k \neq i]} = 1$ iff $k \neq i$. $2n-2$ is the number of samples to contrast in a training batch. As presented in Algorithm 1, we contrast every two sequences in a batch of size n (note that the number of positive pairs is $n-1$), and the final contrastive loss \mathcal{L}_C is computed among all positive pairs.

3.3.2 Analysis on Locality-Aware Contrastive Learning

The proposed LCL aims to improve the learned gait representations AGEs. It is noted that AGEs are constructed with dynamic context vectors of the model, and such context vectors actually play a important role in carrying out the reverse reconstruction of skeleton sequences. Hence, we can also visualize the final reconstruction loss before and after the LCL scheme is applied to improve the learned context vectors, so as to check whether it can improve context vectors and enable better reverse reconstruction. We compare two cases where LCL is applied and not applied (“No LCL”). As shown in Fig. 6, it is found that the LCL constantly enables the model to achieve lower reconstruction loss, regardless of the used attention mechanism (no attention mechanism, BAS,

TABLE 1

Statistics of different datasets. Note: Here we present the number of skeletons estimated from the original CASIA B dataset.

	KGBD	BIWI	KS20	IAS-Lab	CASIA B
# Unique Subjects	164	50	20	11	124
# Original Sequences	822	50	300	11	13639
# Sequences/Classes	5	1	15	1	110
# Skeletons/Classes	2898	369	537	690	8872

MBAS, LAS). These results indicate that the gait encoding model with the LCL scheme can achieve better skeleton reconstruction, and we will empirically demonstrate that gait features learned by LCL can boost the person Re-ID performance as well in Sec. 5.1.

3.3.3 Contrastive Attention-based Gait Encodings

By applying the LCL scheme to sequence-level AGEs (Eq. 12), we can incorporate the inter-sequence locality and tune AGEs into the final gait representations named Contrastive Attention-based Gait Encodings (CAGEs), which preserve locality by both locality-aware attention mechanism and LCL scheme:

$$\mathcal{S}^{(i)} \xrightarrow{\text{I}} \text{AGEs}(\mathbf{V}^{(i)}) \xrightarrow{\text{II}} \text{CAGEs}(\overline{\mathbf{V}}^{(i)}) \quad (15)$$

where $\overline{\mathbf{V}}^{(i)}$ denotes CAGEs of the i^{th} skeleton sequence $\mathcal{S}^{(i)}$, I and II are the learning process of locality-aware attention mechanism and LCL scheme respectively. In this work, our model performs I and II simultaneously at each training step, while we conduct only I in [34]. Note that CAGEs are also built by dynamic context vectors like AGEs in Eq. 11. For simplicity, we use $\overline{\mathbf{V}}$ to represent sequence-level CAGEs, and use \overline{v}_t to represent the skeleton-level CAGE at the t^{th} step of decoding here.

To perform person Re-ID, we use CAGEs to train a simple recognition network $f_{RN}(\cdot)$ that consists of a hidden layer and a softmax layer. In particular, we explore two specific Re-ID strategies: (1) **Sequence-level Concatenation (SC)**: It directly uses sequence-level CAGEs $\overline{\mathbf{V}}$, which is the concatenation of skeleton-level CAGEs ($\{\overline{v}_1; \overline{v}_2; \dots; \overline{v}_f\}$), to train the recognition network and predict the sequence label $f_{RN}(\overline{\mathbf{V}}; \theta_r)$, where θ_r refers to parameters of $f_{RN}(\cdot)$. (2) **Average Prediction (AP)**: It exploits skeleton-level CAGEs to train the recognition network, and averages the prediction of each skeleton-level CAGE $f_{RN}(\overline{v}_t; \theta_r)$ ($t \in \{1, \dots, f\}$) in a skeleton sequence to be the final sequence-level prediction for person Re-ID. We compare AP and SC under different pretext tasks and demonstrate that AP constantly achieves better Re-ID performance than SC (see Sec. 5.2.1). Note that during training, each skeleton in one sequence shares the same skeleton sequence label y_i . Besides, skeleton labels are only used to train the recognition network, *i.e.*, CAGEs are frozen during training. Our later evaluations show that CAGEs, which are learned with unlabeled 3D skeleton data only, are surprisingly discriminative and produce remarkable person Re-ID performance.

3.4 The Entire Approach

As a summary, the computation flow of the entire approach during self-supervised learning is $\mathbf{h} \rightarrow \hat{\mathbf{h}} \rightarrow \overline{\mathbf{a}} \rightarrow \mathbf{c} \rightarrow \overline{\mathbf{h}} \rightarrow \overline{\mathbf{S}}$. To guide model training in the gait encoding process, we combine the loss for self-supervision \mathcal{L}_S (Eq. 3), LA alignment loss \mathcal{L}_A (Eq. 7), and contrastive loss \mathcal{L}_C (Algorithm 1) as follows:

$$\mathcal{L} = \lambda_S \mathcal{L}_S + \lambda_A \mathcal{L}_A + \lambda_C \mathcal{L}_C + \beta \|\Theta\|_2^2 \quad (16)$$

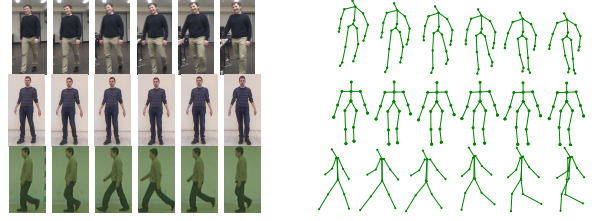


Fig. 7. Examples of RGB images and 3D skeletons in BIWI (first row), IAS-Lab (second row) and CASIA B (third row). Note that the last skeleton sequence is estimated from RGB images of CASIA B.

where Θ denotes the parameters of the model, λ_S , λ_A , λ_C are weight coefficients to trade off the importance of the loss for self-supervision, LA alignment loss and contrastive loss. $\|\Theta\|_2^2$ is L_2 regularization. For the person Re-ID task, we employ standard cross-entropy loss to train the recognition network with CAGEs.

4 EXPERIMENTS

4.1 Experimental Setup

We evaluate our method on four person Re-ID datasets that provide 3D skeleton data: *BIWI* [37], *IAS-Lab* [60], *KS20 VisLab Multi-View Kinect Skeleton Dataset* [35], *Kinect Gait Biometry Dataset (KGBD)* [15], and a large RGB video based multi-view gait dataset *CASIA B* [36]. They collect skeleton data from 50, 11, 20, 164 and 124 different individuals respectively (detailed in Table 3.3.3). As to the former four skeleton-based Re-ID datasets, we follow the evaluation setup in [32], which is frequently used in the literature: For BIWI, we use the training set and *walking* testing set, which contain dynamic skeleton data; For IAS-Lab, we use the full training set and two test splits, IAS-A and IAS-B; For KGBD, since no training and testing splits are given, we randomly leave one skeleton video of each person for testing and use the rest of videos for training. For KS20, we design different split setup to evaluate the multi-view Re-ID performance of our approach: (1) **Random Splits (RS)**: For each viewpoint, we randomly select two skeleton videos for training and use the rest of videos for testing. (2) **Cross-View Splits (CVS)**: We test each viewpoint in KS20 (including left lateral at 0° , left diagonal at 30° , frontal at 90° , right diagonal at 130° , and right lateral at 180°) and use the remaining four viewpoints for training. For each original skeleton sequence that corresponds to an individual person in the dataset, we discard the first and last 10 skeleton frames to avoid ineffective skeleton recording. Then, we split the given original skeleton sequences in the dataset into multiple shorter skeleton sequences (*i.e.*, $\mathcal{S}^{(i)}$) with length f by a step of $\frac{f}{2}$, which aims to obtain as many 3D skeleton sequences as possible to train our model. Unless explicitly specified, the skeleton sequence $\mathcal{S}^{(i)}$ in this paper refers to those split sequences used in learning, rather than those original skeleton sequences provided by datasets.

Different from the aforementioned skeleton-based datasets, CASIA B dataset is a large-scale RGB video based dataset without providing original skeleton data. To apply our method to RGB-based datasets for person Re-ID tasks, we exploit pre-trained pose estimation model [63], [64] to extract 3D skeletons (see Fig. 7) from RGB videos in CASIA B (detailed in supplementary material). CASIA B contains 124 individuals with 11 views— 0° , 18° , 36° , 54° , 72° , 90° , 108° , 126° , 144° , 162° , 180° and three conditions—pedestrians wearing a bag (“bag”), wearing a coat (“clothes”), and without any coat or bag (“normal”). We adopt

TABLE 2

Comparison with existing skeleton-based methods (11-16). Depth-based methods (1-4) and multi-modal methods (5-10) are also included as a reference. Bold numbers refer to the best performers among skeleton-based methods. “—” indicates no published result. “Rev. Rec.” (17) denotes using CAGES learned by the proposed reverse reconstruction, and “Rev. Rec. Plus” (18) represents the proposed enhanced model that concatenates CAGES learned from three pretext tasks for person Re-ID. Best results using average prediction (AP) are reported in 17-18.

		<i>Rank-1</i> (%)					<i>nAUC</i>				
	Id Methods	BIWI	IAS-A	IAS-B	KS20	KGBD	BIWI	IAS-A	IAS-B	KS20	KGBD
Depth-Based	1 Gait Energy Image [20]	21.4	25.6	15.9	—	—	73.2	72.1	66.0	—	—
	2 Gait Energy Volume [21]	25.7	20.4	13.7	—	—	83.2	66.2	64.8	—	—
	3 3D LSTM [32]	27.0	31.0	33.8	—	—	83.3	77.6	78.0	—	—
	4 3D CNN + Average Pooling [61]	27.8	33.4	39.1	—	—	84.0	81.4	82.8	—	—
Multi-Modal	5 PCM + Skeleton [39]	42.9	27.3	81.8	—	—	—	—	—	—	—
	6 Size-Shape descriptors + SVM [42]	20.5	—	—	—	—	—	—	—	—	—
	7 Size-Shape descriptors + LDA [42]	22.1	—	—	—	—	—	—	—	—	—
	8 DVCov + SKL [43]	21.4	46.6	45.9	—	—	—	—	—	—	—
	9 ED + SKL [43]	30.0	52.3	63.3	—	—	—	—	—	—	—
	10 CNN-LSTM with RTA [44]	50.0	—	—	—	—	—	—	—	—	—
Skeleton-Based	11 D^{13} descriptors + SVM [37]	17.9	—	—	—	—	—	—	—	—	—
	12 D^{13} descriptors + KNN [37]	39.3	33.8	40.5	58.3	46.9	64.3	63.6	71.1	78.0	90.0
	13 D^{16} descriptors + Adaboost [38]	41.8	27.4	39.2	59.8	69.9	74.1	65.5	78.2	78.8	90.6
	14 Single-layer LSTM [32]	15.8	20.0	19.1	80.9	39.8	65.8	65.9	68.4	92.3	87.2
	15 Multi-layer LSTM [62]	36.1	34.4	30.9	81.6	46.2	75.6	72.1	71.9	94.2	89.8
	16 PoseGait [14]	33.3	41.4	37.1	70.5	90.6	81.8	79.9	74.8	94.0	97.8
	17 Ours (Rev. Rec.)	62.9	60.1	62.5	86.9	86.9	86.8	82.9	86.9	94.9	97.1
	18 Ours (Rev. Rec. Plus)	63.3	59.1	62.2	92.0	90.6	88.3	81.5	86.2	94.9	98.1

TABLE 3

Re-ID performance comparison on cross-view splits (CVS) of KS20 dataset. 0° , 30° , 90° , 130° , and 180° represent different viewpoints.

	Methods	0°	30°	90°	130°	180°
<i>Rank-1</i>	PoseGait	24.6	19.1	29.7	27.3	25.0
	Ours (Rev. Rec.)	44.4	54.9	55.0	41.9	53.4
	Ours (Rev. Rec. Plus)	48.8	53.6	54.9	44.5	57.5
<i>nAUC</i>	PoseGait	81.2	75.4	81.0	79.6	85.1
	Ours (Rev. Rec.)	84.8	89.1	87.8	83.3	89.3
	Ours (Rev. Rec. Plus)	86.5	87.1	84.8	83.8	91.8

two setups for performance evaluation: (1) Cross-View Evaluation (CVE): We evaluate each view in CASIA B while adjacent views are used for training. (2) Condition-based Matching Evaluation (CME) [65]: We randomly and equally divide 124 people IDs to training set and testing set, and divide the testing set by three original conditions (“normal”, “bag”, “clothes”) to be gallery or probe sets. In particular, we evaluate our approach on single-condition (*i.e.*, gallery set and probe set keep the same condition without appearance changes) and on cross-condition settings (*i.e.*, probe set is under normal condition (“Nm”) while gallery sets are under bag (“Bg”) or clothes condition (“Cl”). Hence, the CME setup contains five combinations, *i.e.*, “Nm-Nm”, “Bg-Bg”, “Cl-Cl”, “Bg-Nm”, and “Cl-Nm”, where the first condition is used for probe set and the second one is for gallery. The details for CME setup can be found in [65]. Experiments under each evaluation setup are repeated for multiple times and the average performance is reported in Sec. 5.4. The implementation details are provided in the supplementary material, and our codes are open at <https://github.com/Kali-Hac/Locality-Awareness-SGE>.

4.2 Evaluation Metrics

Person Re-ID is typically evaluated in a multi-shot manner, and the sequence label can be produced by either predictions of multiple frames or a sequence-level representation. In this work, we report the performance of both strategies (see Table 2 and Table 5): (1) (SC) Using sequence-level CAGES for person Re-ID. (2) (AP) Averaging the prediction of each skeleton-level CAGE in a skeleton sequence to be the final sequence-level prediction for person Re-ID. We compute *Rank-1* accuracy and *nAUC* (area under the cumulative matching curve (CMC) normalized by the number of ranks [66]) to quantify multi-shot person Re-ID performance. It should be noted that we adopt AP strategy for Cross-View Evaluation (CVE) on CASIA B. For Condition-based Matching Evaluation (CME), each sequence-level CAGES in probe set is used to match the one of the same identity in gallery set using Euclidean distance, and the *Rank-1* matching rate is computed.

4.3 Performance Comparison

In this section, we conduct a comprehensive comparison with existing skeleton based person Re-ID methods (Id = 11-16) in the literature. In the meantime, we also include classic depth-based methods (Id = 1-4) and representative multi-modal methods (Id = 5-10) as a reference. The results are reported as follows.

4.3.1 Comparison with Skeleton-based Methods

As shown by Table 2, our approach enjoys obvious advantages over existing skeleton-based methods in terms of both Re-ID performance metrics: First, our approach evidently outperforms those methods that rely on manually-designed geometric or anthropometric skeleton descriptors (Id = 11-13). For example, D^{13} (Id = 12) and recent D^{16} (Id = 13) are two most representative hand-crafted feature based methods, and our model outperforms both of them by a large margin (20.7%-43.7% *Rank-1* accuracy

TABLE 4

Ablation study of our model. “✓” indicates that the corresponding model component is used: GE, GD, reverse skeleton reconstruction (Rev. Rec.), locality-aware attention alignment scores (LAS). “AGEs” indicates using AGEs (v_t) rather than encoded gait states of GE’s LSTM h_t to perform person Re-ID task (note that here h_t and AGEs (v_t) are learned without the locality-aware contrastive learning (LCL) scheme), and “LCL+CAGEs” (\bar{v}_t) represents exploiting CAGEs learned by the LCL scheme for person Re-ID.

Model Configuration						BIWI		IAS-A		IAS-B		KS20	
GE	GD	Rev. Rec.	LAS	AGEs	LCL+CAGEs	Rank-1	nAUC	Rank-1	nAUC	Rank-1	nAUC	Rank-1	nAUC
✓						36.1	75.6	34.4	72.1	30.9	71.9	80.9	92.3
✓	✓					41.5	80.1	48.1	77.5	48.4	76.2	83.3	92.2
✓	✓	✓				46.7	81.5	50.9	78.3	52.9	80.3	84.5	94.2
✓	✓	✓	✓			57.7	85.8	55.4	81.6	57.4	83.6	85.7	94.1
✓	✓		✓	✓		57.2	85.7	55.6	80.7	57.0	84.8	86.5	94.7
✓	✓	✓	✓	✓		59.1	86.5	56.1	80.7	58.2	85.3	86.7	93.2
✓	✓		✓		✓	59.7	86.6	57.9	82.7	60.9	84.3	85.9	94.7
✓	✓	✓	✓		✓	62.9	86.8	60.1	82.9	62.5	86.2	86.9	95.7

and 7.5%-24.0% *nAUC* on different datasets). Second, our approach is also superior to recent skeleton based methods that utilize deep neural networks (Id = 14-16) on all datasets by up to 50.8% *Rank-1* accuracy and 22.5% *nAUC* improvement. Although the latest PoseGait can achieve comparable performance to our approach on the KGBD dataset, it still requires extracting 81 hand-crafted features for CNN learning. Besides, labeled skeleton data are indispensable for the gait encoding stage of existing deep learning based methods, while our approach can learn better gait representations by simply exploiting unlabeled 3D skeleton data.

Besides, we also evaluate the cross-view Re-ID performance of our approach on the KS20 dataset using multi-view 3D skeleton data provided by this dataset. We compare its performance with the latest PoseGait [14] approach, which specially considers the multi-view scenario and achieves the best overall performance among existing skeleton based methods. The results are displayed in Table 3 and highlight the following conclusions: Our approach consistently outperforms PoseGait by a considerable margin (up to 35.8% *Rank-1* accuracy and 13.7% *nAUC*) under all viewpoints. In the meantime, it can stably achieve comparatively satisfactory performance under different viewpoints, which validates the robustness of our approach against viewpoint variations.

4.3.2 Comparison with Depth-based Methods and Multi-modal Methods

Despite that our approach only takes 3D skeleton data as inputs, our approach consistently outperforms baselines of classic depth-based methods (Id = 1-4) by at least 23.4% *Rank-1* and 1.5% *nAUC* gain. Considering the fact that 3D skeletons are of much smaller data size than depth image data, our approach is both effective and efficient. As to the comparison with recent methods that exploit multi-modal inputs (Id = 5-10), the performance of our approach is still highly competitive: Although in few cases multi-modal methods perform better on IAS-B, our skeleton based method achieves the best *Rank-1* accuracy on BIWI and IAS-A. Interestingly, we note that the multi-modal approach that uses both point cloud matching (PCM) and skeletons yields the best accuracy on IAS-B, but it performs markedly worse on datasets that undergo more frequent shape and appearance changes (IAS-A and BIWI). By contrast, our approach consistently achieves stable and satisfactory performance on each dataset. Thus, with 3D skeleton data as the sole input, our approach can be a promising solution to person Re-ID and other potential skeleton-related tasks.

5 FURTHER ANALYSIS

5.1 Ablation Studies

In this section, we carry out ablation studies to verify the necessity of each model component in the proposed approach. As shown in Table 4, we can arrive at the following conclusions: **(1)** The proposed encoder-decoder architecture (GE-GD) performs remarkably better than the supervised learning paradigm, which uses GE only to perform person Re-ID directly (up to 17.5% *Rank-1* accuracy and 5.4% *nAUC* gain). Such results demonstrate the effectiveness of our self-supervised gait encoding model, which leverages an encoder-decoder architecture and skeleton sequence reconstruction mechanism. **(2)** Using reverse reconstruction (“Rev. Rec.” in Table 4) typically produces evident performance gain (up to 5.2% *Rank-1* accuracy and 4.1% *nAUC*) when compared with those configurations without reverse reconstruction. Such results justify reverse reconstruction as an effective pretext task for gait encoding, as it enables the model to learn more discriminative gait features for person Re-ID. **(3)** Adding the proposed locality-aware attention mechanism (“LAS”) is able to improve the performance remarkably by up to 11.0% *Rank-1* accuracy and 4.3% *nAUC*. This observation is consistent with our previous analysis that LAS facilitates reverse reconstruction and contributes to person Re-ID performance. We will compare different attention mechanisms in the next section as well. **(4)** The proposed context vector based gait representations, CAGEs and AGEs, are both able to achieve superior performance to frequently-used features (h_t). When it comes to the comparison between AGEs (learned without contrastive learning) and CAGEs (learned with locality-aware contrastive learning), CAGEs consistently outperforms AGEs with a 0.2%-4.3% *Rank-1* accuracy and 0.3%-2.5% *nAUC* gain on different datasets (see row 6 and row 8 in Table 4). This demonstrates the advantage to incorporate inter-sequence locality by the proposed locality-aware contrastive learning scheme.

5.2 Discussions

5.2.1 Different Pretext Tasks

In this section, we systematically explore the influence of pretext tasks on self-supervised learning in Table 5. First, we evaluate the Re-ID performance of different pretext tasks (Prediction, Sorting, and Rev. Rec.) under two Re-ID strategies (average prediction (AP) and sequence-level concatenation (SC)) on different datasets. Note that we apply the basic attention mechanism to prediction

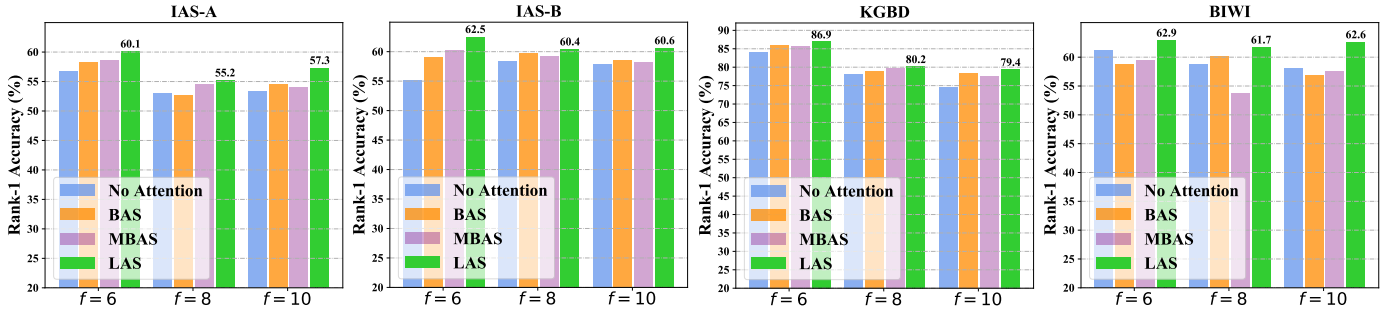


Fig. 8. Rank-1 accuracy on different datasets when using no attention, BAS, MBAS or LAS for model learning. f denotes the sequence length.

TABLE 5

Performance comparison of different pretext tasks and the proposed enhanced approach under two Re-ID manners (“AP”: Average prediction. “SC”: Sequence-level concatenation). Bold numbers refer to the best Rank-1 accuracy and $nAUC$ among different configurations.

Dataset	Pretext Task(s)	Rank-1		$nAUC$	
		AP	SC	AP	SC
BIWI	Prediction	40.4	32.0	82.0	79.7
	Sorting	55.7	43.4	71.4	85.5
	Rev. Rec.	62.9	51.3	86.8	84.1
	Pred. + Sort.	62.1	51.5	86.7	86.9
	Rev. Rec. + Pred.	62.3	52.5	87.5	86.6
	Rev. Rec. + Sort.	63.1	51.3	88.0	86.0
	Rev. Rec. Plus	63.3	53.7	88.3	87.1
IAS-A	Prediction	56.7	54.6	82.3	72.9
	Sorting	56.8	51.9	82.9	78.9
	Rev. Rec.	60.1	55.0	82.9	81.1
	Pred. + Sort.	58.8	54.0	80.6	80.1
	Rev. Rec. + Pred.	58.8	52.6	81.1	79.1
	Rev. Rec. + Sort.	58.7	54.1	80.9	80.3
	Rev. Rec. Plus	59.1	53.0	81.5	79.3
IAS-B	Prediction	58.5	55.8	84.6	76.0
	Sorting	54.1	49.5	83.6	78.5
	Rev. Rec.	62.5	53.5	86.9	85.6
	Pred. + Sort.	61.9	54.7	86.2	84.7
	Rev. Rec. + Pred.	61.2	56.7	86.7	85.7
	Rev. Rec. + Sort.	61.5	55.0	86.0	85.1
	Rev. Rec. Plus	62.2	57.6	86.2	85.2
KS20	Prediction	81.6	74.0	77.5	89.9
	Sorting	83.5	78.1	90.6	91.7
	Rev. Rec.	86.9	84.1	94.9	85.7
	Pred. + Sort.	89.5	86.3	92.8	88.4
	Rev. Rec. + Pred.	91.9	86.0	94.5	90.1
	Rev. Rec. + Sort.	90.9	84.7	93.1	87.4
	Rev. Rec. Plus	92.0	89.4	94.9	89.5
KGBD	Prediction	85.5	85.1	96.5	97.9
	Sorting	85.4	84.2	97.4	98.0
	Rev. Rec.	86.9	84.1	97.1	97.4
	Pred. + Sort.	90.4	87.6	97.9	97.8
	Rev. Rec. + Pred.	90.6	88.5	97.1	97.7
	Rev. Rec. + Sort.	90.5	87.5	97.4	96.9
	Rev. Rec. Plus	90.6	88.3	98.1	98.0

and sorting, because they cannot exploit intra-sequence locality and using locality-aware attention actually does not improve their performance. Then, we also compare their performance with all potential combinations of pretext tasks: (1) Pred. + Sort. , (2) Rev. Rec. + Pred. , (3) Rev. Rec. + Sort. and (4) the proposed

enhanced configuration Rev. Rec. Plus (*i.e.*, Rev. Rec. + Pred. + Sort.) (introduced in Sec. 3.1.2). The results are exhibited in Table 5, and we draw the following conclusions: **(1)** When a single pretext task is used for self-supervised learning, reverse reconstruction typically performs comparably or superior to other pretext tasks in terms of both AP and SC. This is due to the fact that reverse reconstruction utilizes the skeleton order information embedded in inputs, which enables the exploitation of the intra-sequence locality during training. Such results also justify the center role of reverse reconstruction in self-supervised learning. **(2)** Combining gait features learned from different pretext tasks could achieve higher Re-ID performance than using Rev. Rec. alone in many cases. Notably, the enhanced configuration Rev. Rec. Plus, which combines CAGEs learned from all three pretext tasks, obtains the best overall performance in terms of Rank-1 and $nAUC$ on three out of four datasets (BIWI, KS20, KGBD). Such observations reveal the potential to extract richer gait features for performance improvement through introducing more pretext tasks into self-supervised learning. Nevertheless, Rev. Rec. Plus requires learning three different pretext tasks, and suffers from higher computational cost and feature dimension. By contrast, reverse reconstruction usually obtains similarly competitive performance with much less training cost and simpler gait representations. Hence, we recommend to use reverse reconstruction as the primary pretext task, and our later analysis is also performed based on this pretext task. **(3)** As shown in Table 5, using AP almost constantly achieves better Re-ID performance. The reason is that AP is able to reduce the influence of noisy frames that give wrong predictions, so it can encourage better sequence-level predictions.

5.2.2 Different Attention Mechanisms

In order to show the effects of attention mechanisms, we evaluate the person Re-ID performance of our approach under four different cases (no attention, BAS, MBAS, or LAS). To provide a more comprehensive evaluation, we also evaluate attention mechanisms with different sequence lengths. By results reported in Fig. 8, we can draw the following conclusions: **(1)** The application of attention mechanisms (BAS, MBAS, LAS) can improve the model performance by 1.0%-7.2% Rank-1 accuracy when compared with the case without attention mechanism. This is because the attention mechanisms can help the model focus on more correlative skeletons, thus leading to better sequence reconstruction and more effective gait representations for person Re-ID. **(2)** Among different attention mechanisms, the proposed locality-aware attention mechanism (LAS) is the best performer, which surpasses BAS and MBAS by up to 7.9% Rank-1 accuracy on different datasets.

TABLE 6

Performance of our approach when setting different contrasting intervals for learning (“Interval=1” indicates contrasting adjacent sequences).

Interval	Rank-1					nAUC				
	BIWI	IAS-A	IAS-B	KS20	KGBD	BIWI	IAS-A	IAS-B	KS20	KGBD
1	62.9	60.1	62.5	86.9	86.9	86.8	82.9	86.9	94.9	97.1
2	61.1	57.9	61.4	86.5	86.4	86.5	82.5	84.8	94.4	96.8
3	61.9	59.4	58.1	86.1	86.7	86.0	82.3	84.8	94.1	96.8
4	60.8	59.2	58.9	85.9	86.5	86.7	82.4	82.7	94.9	96.7

TABLE 7

Performance of our approach when setting different temperatures ($\tau = 0.05, 0.1, 0.5, 0.8, 1$) for the LCL scheme on different datasets.

τ	Rank-1					nAUC				
	BIWI	IAS-A	IAS-B	KS20	KGBD	BIWI	IAS-A	IAS-B	KS20	KGBD
1	60.6	58.6	59.7	85.9	86.8	85.3	81.3	84.5	94.5	97.1
0.8	61.2	59.3	60.2	86.3	86.6	85.6	82.3	85.2	95.0	97.0
0.5	61.1	58.5	61.5	86.3	86.9	86.8	82.0	85.8	95.0	97.1
0.1	62.9	60.1	62.5	86.9	86.8	86.8	82.9	86.9	94.9	96.6
0.05	62.6	59.0	61.9	86.7	86.6	86.3	83.0	86.4	95.2	96.8

These results justify our claim that intra-sequence locality enables better gait representation learning for person Re-ID.

5.2.3 Different Contrasting Intervals

In this section, we discuss the performance of our approach under different contrasting intervals. To be more specific, here we not only contrast adjacent skeleton sequences (Interval=1) by the proposed LCL scheme, but also exploit those non-adjacent skeleton sequences (Interval>1) for contrastive learning. As shown in Table 6, when compared with other contrasting interval settings, using adjacent sequences to perform LCL scheme constantly achieves the best *Rank-1* accuracy as well as *nAUC* on all datasets. The comparison validates that adjacent skeleton sequences in a local temporal context enjoy higher correlations. In other words, such results justify our motivation to integrate the inter-sequence locality, which can be effectively learned by our LCL scheme, for the enhancement of the gait encoding for person Re-ID.

5.2.4 Temperature Setting for LCL Scheme

When the proposed LCL scheme is applied, we need to determine the value of temperature τ . In this section, we evaluate the performance of our approach when different values are set to the temperature τ of LCL scheme. As shown in Table 7, we observe that the performance of our model typically enjoys a stable performance when the value of τ is varied, and no drastic change is observed. The results suggest that our LCL scheme is actually insensitive to τ , so we simply set τ empirically in our experiments.

5.3 Transferability of Gait Encoding Model

Interestingly, we discover that the gait encoding model learned on one dataset can be readily transferred to other datasets. Specifically, we use the gait encoding model pre-trained on training sets (“source datasets”) from KGBD, BIWI or IAS-Lab to directly encode 3D skeleton data from other datasets (“target datasets”) into CAGEs, which are then used for training the recognition network f_{RN} to perform person Re-ID. We compare their Re-ID performance with the gait encoding model trained on their original

training sets: As shown by Fig. 9, the transferred gait encoding models (*i.e.*, trained on a different source dataset) are also able to achieve highly competitive performance when compared with existing methods in literature, despite that our original model is still the best performer. Such transferability enables the pre-trained model to extract discriminative gait features from unseen skeleton data of a new dataset, which demonstrates that our approach indeed captures transferable high-level semantics of 3D skeleton data.

5.4 Evaluation on Model-Estimated Skeleton Data

To further evaluate our skeleton-based approach on the large-scale gait dataset CASIA B, we exploit pre-trained pose estimation models [63], [64] to extract skeleton data from RGB videos of CASIA B, and evaluate the performance of our approach with the estimated skeleton data. We compare our model with the latest skeleton-based method PoseGait [14] and representative appearance-based methods [65]–[69], and we can draw the following observations and conclusions:

(1) As shown in Table 8, our approach outperforms the latest skeleton-based model PoseGait by a significant margin (24.8%–57.1% *Rank-1* accuracy) on all views of CASIA B. Notably, both methods achieve their own best performance on the view 36°, while PoseGait is still inferior to our method by 27% *Rank-1* accuracy. On the two most challenging views (0° and 180°), our approach also achieves better Re-ID performance than PoseGait by more than 24% *Rank-1* accuracy. (2) In Table 9, our skeleton-based approach could also outperform representative classic appearance-based methods that utilize visual features (*e.g.*, RGB features, silhouettes). For example, our method outperforms LMNN [67] and ITML [68], which leverage visual features (RGB, HSV color and texture) and metric learning for recognition [65], with an evident performance gain up to 50.4% on different CME settings. Compared with the method that fuses RGB appearance features and GEI features (“Feature-level + MLR” [65]), our model also yields evidently better *Rank-1* matching rate by 4.5%–38.0% on all conditions of CASIA B. Therefore, despite that our

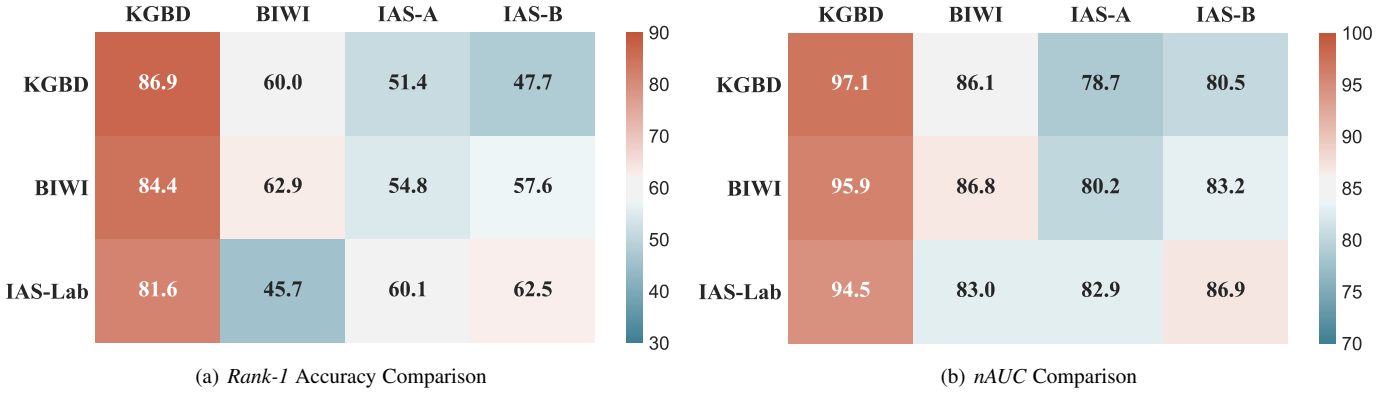


Fig. 9. Rank-1 accuracy and $nAUC$ comparison between the original model and the transferred model on different datasets. Note that the abscissa and ordinate denote target datasets and source datasets (for training gait encoding models) respectively.

TABLE 8
Rank-1 accuracy on different views of CASIA B under CVE setup.

Methods	0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°
PoseGait [14]	10.7	37.4	52.5	28.3	24.3	18.9	23.5	17.2	23.6	18.8	4.3
Ours	35.5	78.5	79.5	58.0	66.1	76.0	64.4	66.9	68.4	49.9	46.3

TABLE 9
Rank-1 matching rate on CASIA B compared with appearance-based methods under CME setup. Note: “Cl-Nm” denotes the probe set (under “Clothes” condition) and gallery set (under “Normal” condition).

Methods	Nm-Nm	Bg-Bg	Cl-Cl	Cl-Nm	Bg-Nm
LMNN [67]	3.9	18.3	17.4	11.6	23.1
ITML [68]	7.5	19.5	20.1	10.3	21.8
EFL [66]	12.3	5.8	19.9	5.6	17.1
SDALF [69]	4.9	10.2	16.7	11.6	22.9
Score-Level + MLR [65]	13.6	13.6	13.5	9.7	14.7
Feature-level + MLR [65]	16.3	18.9	25.4	20.3	31.8
Ours	54.3	37.5	31.9	27.0	36.3

approach is trained on the estimated skeleton data with noises, the learned gait representations still show higher discriminative power and achieve superior performance to those appearance-based methods on CASIA B. These results show the effectiveness of our approach on multi-view Re-ID tasks when using model-estimated skeleton data, and also demonstrate the great potential of our approach to be applied to large RGB-based datasets.

5.5 Visualization of Typical Samples

In Fig. 11, we visualize both simple samples and hard samples for Re-ID from testing sequences of different datasets. To this end, we select the classes that are most likely to be confused (*e.g.*, ID 6 and 7 in BIWI, shown in Fig. 10 (b)) for visualization. We obtain the following observations: (1) Some skeletons from those datasets indeed contain noise. For example, in Fig. 11 (b), the first skeleton of the sequence suffers from obvious twisting, which can be ascribed to wrongly detected body joints, while similar noise can be observed from the former three skeletons in Fig. 11 (f). In the mean time, we notice that skeleton sequences contaminated by noise are often wrongly classified, which suggest that noise degrades the person Re-ID performance. (2) The skeleton sequences that contain salient motion are more likely to be

correctly classified. For example, a skeleton sequence with intense action (see Fig. 11 (c)) can be easily classified by our model, while those sequences with comparatively static skeletons are often hard to recognize (see Fig. 11 (d)). This actually validates that salient gait patterns in skeleton sequences will facilitate person Re-ID.

6 CONCLUSION AND FUTURE WORK

In this paper, we propose a generic self-supervised approach with locality-awareness to learn effective gait representations for person Re-ID. We introduce self-supervision by learning reverse skeleton sequence reconstruction as a primary pretext task, which enables our model to learn high-level semantics and discriminative gait features with unlabeled skeleton data. Other potential pretext tasks like sorting and prediction are also explored and synthesized into the self-supervised learning. To facilitate skeleton reconstruction and gait representation learning, a novel locality-aware attention mechanism and locality-aware contrastive learning scheme are proposed to incorporate the intra-sequence and inter-sequence locality into gait encoding process. Last, we propose to construct the final gait representations (CAGEs) for person Re-ID with learned context vectors. Our approach significantly outperforms existing skeleton-based Re-ID methods, and its performance is comparable or superior to depth-based and multi-modal methods. Besides, we show that our approach can be applied to 3D skeleton data estimated from large RGB-based datasets, and achieve better performance than many classic appearance-based methods.

Limitations. There are three limitations in our work: First, the scale of skeleton datasets in our experiments is relatively limited when compared with popular RGB-based Re-ID datasets like DukeMTMC-reID, since large-scale skeleton-based Re-ID datasets that contain more individuals and scenarios are still unavailable. Hopefully, we will collect our own skeleton-based Re-ID datasets in the future. Second, this work basically considers the case where the skeletons are collected with relatively high quality (*e.g.*, by device like Kinect), while the case where skeleton data are collected under a more general setting (*e.g.*, estimated

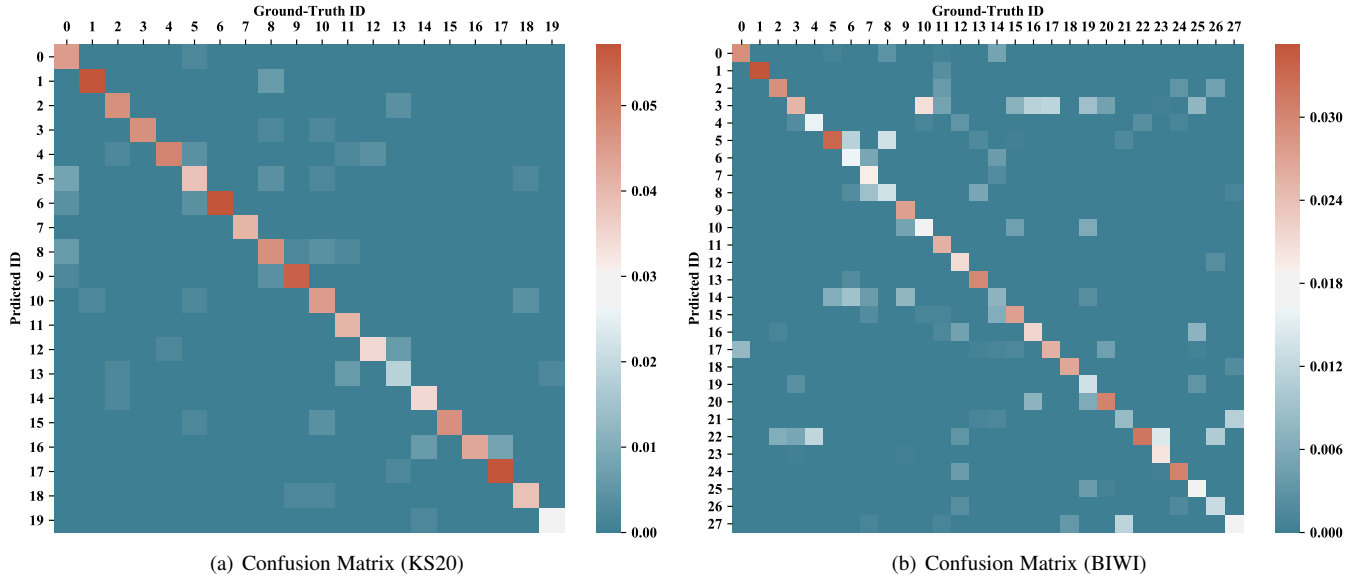


Fig. 10. Confusion matrices of KS20 and BIWI. Note that abscissa and ordinate denote the ground-truth and predicted IDs respectively. The position in the a^{th} column and b^{th} row indicates that the testing samples belonging to the a^{th} ID is predicted as the b^{th} ID, while the corresponding value is the proportion of such samples to samples in the whole testing set. Full results are provided in the supplementary material.

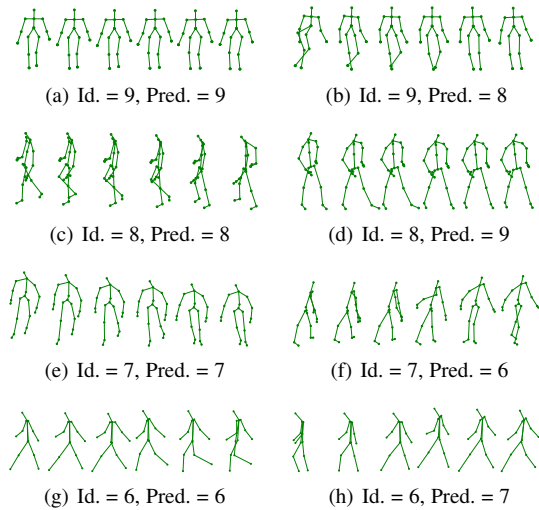


Fig. 11. Visualization of typical samples in testing sets of IAS-B (a-b), KS20 (c-d), BIWI (e-f), CASIA B (g-h). Note: “Id.” and “Pred.” denotes the ground-truth label and predicted label respectively. (g) and (h) are partial samples of testing sequences on the view 90° of CASIA B.

from RGB data in outdoor scenes) has not been thoroughly studied. Finally, our approach models 3D skeletons as body-joint sequences on different dimensions, while it may be insufficient to capture underlying relations between body joints.

Overall, our approach showcases the effectiveness of self-supervised gait encoding on the person Re-ID task, and there are several potential directions for improvement: **(1)** More efficient pretext tasks (*e.g.*, frame interpolation, skeleton video generation) could be explored to improve the capture of motion semantics for better gait encoding. **(2)** Modeling 3D skeletons as graphs is able to mine richer relation information among body joints, while employing graph-based encoders (*e.g.*, graph convolutional network (GCN)) could enhance structural feature learning from skeleton graphs. **(3)** Fine-grained spatial-temporal attention mechanisms

could also be designed to extract those crucial motion patterns for person Re-ID, and more effective skeleton augmentation strategies could be considered to enhance the contrastive learning. **(4)** One important future direction is to explore 3D skeleton-based Re-ID under the general setting, so as to improve the model robustness to noise, *e.g.*, noise incurred by skeleton/pose estimation. **(5)** Our model can be extended to more skeleton-related tasks, and we can expect it to be readily transferred to multi-modal learning for other pivotal vision tasks.

7 ETHICAL STATEMENTS

Person Re-ID is an important topic with huge potential value in computer vision. However, it should be noted that improper application or abuse of Re-ID technology will pose a grave threat to the society and public privacy. Thus, we want to emphasize that benchmark datasets used in this work are either publicly available (BIWI, IAS-Lab, KGBD) or officially authorized (KS20, CASIA B). The official agents of those datasets have guaranteed that all data are collected, released, and used with the consent of subjects. All people in datasets are anonymous with simple identity numbers for privacy protection. Besides, our approach and models must only be used for research purpose.

ACKNOWLEDGMENTS

This work was supported in part by the National Key Research and Development Program of China (Grant No. 2019YFA0706200), in part by the National Natural Science Foundation of China (Grant No. 61632014, No. 61627808, No. 62006236, No. 62072190), in part by the Hunan Provincial Natural Science Foundation (Grant No. 2020JJ5673), and in part by the NUDT Research Project (Grant No. ZK20-10).

REFERENCES

- [1] C. C. Loy, T. Xiang, and S. Gong, “Time-delayed correlation analysis for multi-camera activity understanding,” *International Journal of Computer Vision*, vol. 90, no. 1, pp. 106–129, 2010.

- [2] D. Baltieri, R. Vezzani, and R. Cucchiara, "Sarc3d: a new 3d body model for people tracking and re-identification," in *International Conference on Image Analysis and Processing*. Springer, 2011, pp. 197–206.
- [3] R. Vezzani, D. Baltieri, and R. Cucchiara, "People reidentification in surveillance and forensics: A survey," *ACM Computing Surveys*, vol. 46, no. 2, p. 29, 2013.
- [4] T. Wang, S. Gong, X. Zhu, and S. Wang, "Person re-identification by discriminative selection in video ranking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 12, pp. 2501–2514, 2016.
- [5] R. Zhao, W. Oyang, and X. Wang, "Person re-identification by saliency learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 2, pp. 356–370, 2017.
- [6] C. Su, F. Yang, S. Zhang, Q. Tian, L. S. Davis, and W. Gao, "Multi-task learning with low rank attribute embedding for multi-camera person re-identification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1167–1181, 2018.
- [7] Y.-C. Chen, X. Zhu, W.-S. Zheng, and J.-H. Lai, "Person re-identification by camera correlation aware feature augmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 2, pp. 392–408, 2018.
- [8] M. Li, X. Zhu, and S. Gong, "Unsupervised tracklet person re-identification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2019.
- [9] X. Qian, Y. Fu, T. Xiang, Y.-G. Jiang, and X. Xue, "Leader-based multi-scale attention deep architecture for person re-identification," *IEEE Transaction on Pattern Analysis Machine Intelligence*, 2019.
- [10] H.-X. Yu, A. Wu, and W.-S. Zheng, "Unsupervised person re-identification by deep asymmetric metric embedding," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 4, pp. 956–973, 2020.
- [11] P. Connor and A. Ross, "Biometric recognition by gait: A survey of modalities and features," *Computer Vision and Image Understanding*, vol. 167, pp. 1–27, 2018.
- [12] M. P. Murray, A. B. Drought, and R. C. Kory, "Walking patterns of normal men," *Journal of Bone and Joint Surgery*, vol. 46, no. 2, pp. 335–360, 1964.
- [13] J. E. Cutting and L. T. Kozlowski, "Recognizing friends by their walk: Gait perception without familiarity cues," *Bulletin of the psychonomic society*, vol. 9, no. 5, pp. 353–356, 1977.
- [14] R. Liao, S. Yu, W. An, and Y. Huang, "A model-based gait recognition method with body pose and human prior knowledge," *Pattern Recognition*, vol. 98, p. 107069, 2020.
- [15] V. O. Andersson and R. M. Araujo, "Person identification using anthropometric and gait data from kinect sensor," in *AAAI*, 2015.
- [16] L. Wang, T. Tan, H. Ning, and W. Hu, "Silhouette analysis-based gait recognition for human identification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 12, pp. 1505–1518, 2003.
- [17] A. Veeraraghavan, A. K. Roy-Chowdhury, and R. Chellappa, "Matching shape sequences in video with applications in human movement analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1896–1909, 2005.
- [18] Z. Liu and S. Sarkar, "Improved gait recognition by gait dynamics normalization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 6, pp. 863–876, 2006.
- [19] Y. Guan, C.-T. Li, and F. Roli, "On reducing the effect of covariate factors in gait recognition: a classifier ensemble method," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 7, pp. 1521–1528, 2014.
- [20] L. Chunli and W. Kejun, "A behavior classification based on enhanced gait energy image," in *International Conference on Networking and Digital Society*, vol. 2. IEEE, 2010, pp. 589–592.
- [21] S. Sivapalan, D. Chen, S. Denman, S. Sridharan, and C. Fookes, "Gait energy volumes and frontal gait recognition using depth images," in *International Joint Conference on Biometrics*. IEEE, 2011, pp. 1–6.
- [22] C. Wang, J. Zhang, L. Wang, J. Pu, and X. Yuan, "Human identification using temporal information preserving gait template," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 11, pp. 2164–2176, 2011.
- [23] Y. Zhang, Y. Huang, L. Wang, and S. Yu, "A comprehensive study on gait biometrics using a joint cnn-based method," *Pattern Recognition*, vol. 93, pp. 228–236, 2019.
- [24] R. Tanawongsuwan and A. Bobick, "Gait recognition from time-normalized joint-angle trajectories in the walking plane," in *CVPR*, vol. 2, Dec 2001, pp. II–II.
- [25] L. Wang, H. Ning, T. Tan, and W. Hu, "Fusion of static and dynamic body biometrics for gait recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 14, no. 2, pp. 149–158, 2004, cited By 237.
- [26] G. Ariyanto and M. S. Nixon, "Model-based 3d gait biometrics," in *International Joint Conference on Biometrics*, Oct 2011, pp. 1–7.
- [27] A. Nambiar, A. Bernardino, and J. C. Nascimento, "Gait-based person re-identification: A survey," *ACM Computing Surveys*, vol. 52, no. 2, p. 33, 2019.
- [28] F. Han, B. Reily, W. Hoff, and H. Zhang, "Space-time representation of people based on 3d skeletal data: A review," *Computer Vision and Image Understanding*, vol. 158, pp. 85–105, 2017.
- [29] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. J. Finocchio, R. Moore, A. A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *CVPR*, 2011, pp. 1297–1304.
- [30] I. B. Barbosa, M. Cristani, A. Del Bue, L. Bazzani, and V. Murino, "Re-identification with rgb-d sensors," in *ECCV*. Springer, 2012, pp. 433–442.
- [31] J.-H. Yoo, M. S. Nixon, and C. J. Harris, "Extracting gait signatures based on anatomical knowledge," in *Proceedings of BMVA Symposium on Advancing Biometric Technologies*. Citeseer, 2002, pp. 596–606.
- [32] A. Haque, A. Alahi, and L. Fei-Fei, "Recurrent attention models for depth-based person identification," in *CVPR*, 2016, pp. 1229–1238.
- [33] J. K. Aggarwal and Q. Cai, "Human motion analysis: A review," *Computer Vision and Image Understanding*, vol. 73, no. 3, pp. 428–440, 1999.
- [34] H. Rao, S. Wang, X. Hu, M. Tan, H. Da, J. Cheng, and B. Hu, "Self-supervised gait encoding with locality-aware attention for person re-identification," in *IJCAI*, C. Bessiere, Ed. ijcai.org, 2020, pp. 898–905.
- [35] A. Nambiar, A. Bernardino, J. C. Nascimento, and A. Fred, "Context-aware person re-identification in the wild via fusion of gait and anthropometric features," in *International Conference on Automatic Face & Gesture Recognition*. IEEE, 2017, pp. 973–980.
- [36] 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*, vol. 4. IEEE, 2006, pp. 441–444.
- [37] M. Munaro, A. Fossati, A. Basso, E. Menegatti, and L. Van Gool, "One-shot person re-identification with a consumer depth camera," in *Person Re-Identification*. Springer, 2014, pp. 161–181.
- [38] P. Pala, L. Seidenari, S. Berretti, and A. Del Bimbo, "Enhanced skeleton and face 3d data for person re-identification from depth cameras," *Computers & Graphics*, 2019.
- [39] M. Munaro, A. Basso, A. Fossati, L. Van Gool, and E. Menegatti, "3d reconstruction of freely moving persons for re-identification with a depth sensor," in *ICRA*. IEEE, 2014, pp. 4512–4519.
- [40] F. Battistone and A. Petrosino, "Tglstm: A time based graph deep learning approach to gait recognition," *Pattern Recognition Letters*, 2018.
- [41] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [42] M. Hasan and N. Babaguchi, "Long-term people reidentification using anthropometric signature," in *International Conference on Biometrics Theory, Applications and Systems*. IEEE, 2016, pp. 1–6.
- [43] A. Wu, W.-S. Zheng, and J.-H. Lai, "Robust depth-based person re-identification," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp. 2588–2603, 2017.
- [44] N. Karianakis, Z. Liu, Y. Chen, and S. Soatto, "Reinforced temporal attention and split-rate transfer for depth-based person re-identification," in *ECCV*, 2018, pp. 715–733.
- [45] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *CVPR*, vol. 2, 2006, pp. 1735–1742.
- [46] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *CVPR*, 2018, pp. 3733–3742.
- [47] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [48] C. Zhuang, A. L. Zhai, and D. Yamins, "Local aggregation for unsupervised learning of visual embeddings," in *ICCV*, 2019, pp. 6002–6012.
- [49] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox, "Discriminative unsupervised feature learning with convolutional neural networks," in *NeurIPS*, 2014, pp. 766–774.
- [50] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *International Conference on Artificial Intelligence and Statistics*, 2010, pp. 297–304.
- [51] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," *arXiv preprint arXiv:1906.05849*, 2019.
- [52] I. Misra and L. van der Maaten, "Self-supervised learning of pretext-invariant representations," in *CVPR*, 2020, pp. 6707–6717.

- [53] M. Ye, X. Zhang, P. C. Yuen, and S.-F. Chang, "Unsupervised embedding learning via invariant and spreading instance feature," in *CVPR*, 2019, pp. 6210–6219.
- [54] X. Ji, J. F. Henriques, and A. Vedaldi, "Invariant information clustering for unsupervised image classification and segmentation," in *ICCV*, 2019, pp. 9865–9874.
- [55] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, 2020.
- [56] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. E. Hinton, "Big self-supervised models are strong semi-supervised learners," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 22243–22255. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/fcbe95ccdd51da181207c0e1400c655-Paper.pdf>
- [57] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020, pp. 9729–9738.
- [58] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *EMNLP*, Sep. 2015, pp. 1412–1421.
- [59] J. Weston, S. Chopra, and A. Bordes, "Memory networks," in *ICLR*, 2015.
- [60] M. Munaro, S. Ghidoni, D. T. Dizmen, and E. Menegatti, "A feature-based approach to people re-identification using skeleton keypoints," in *ICRA*. IEEE, 2014, pp. 5644–5651.
- [61] Y.-L. Boureau, J. Ponce, and Y. LeCun, "A theoretical analysis of feature pooling in visual recognition," in *ICML*, 2010, pp. 111–118.
- [62] W. Zheng, L. Li, Z. Zhang, Y. Huang, and L. Wang, "Relational network for skeleton-based action recognition," in *ICME*. IEEE, 2019, pp. 826–831.
- [63] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "Openpose: realtime multi-person 2d pose estimation using part affinity fields," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 1, pp. 172–186, 2019.
- [64] C.-H. Chen and D. Ramanan, "3d human pose estimation= 2d pose estimation+ matching," in *CVPR*, 2017, pp. 7035–7043.
- [65] Z. Liu, Z. Zhang, Q. Wu, and Y. Wang, "Enhancing person re-identification by integrating gait biometric," *Neurocomputing*, vol. 168, pp. 1144–1156, 2015.
- [66] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an ensemble of localized features," in *ECCV*. Springer, 2008, pp. 262–275.
- [67] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification." *Journal of machine learning research*, vol. 10, no. 2, 2009.
- [68] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *ICML*, 2007, pp. 209–216.
- [69] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani, "Person re-identification by symmetry-driven accumulation of local features," in *CVPR*. IEEE, 2010, pp. 2360–2367.



Haocong Rao received the B.Eng degree from South China University of Technology, China, in 2019. He was a visiting student with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China, in 2020. He is currently working toward the PhD degree at Nanyang Technological University, Singapore. He has authored/coauthored more than 5 peer-reviewed papers in highly regarded journals and conferences such as IJCAI, ACMMM, IEEE IoT journal, Information Sciences, etc. He received

the Best Paper Award from IEEE Heathcom 2020. His research interests include skeleton-based person re-identification, self-supervised learning, domain adaptation and interpretable artificial intelligence.



Siqi Wang received the BS degree and the PhD degree in computer science and technology from the National University of Defense Technology, China. He is currently an assistant professor with the State Key Laboratory of High Performance Computing (HPCL), National University of Defense Technology, China. His main research focuses on anomaly / outlier detection, pattern recognition and unsupervised learning. His works have been published on leading conferences and journals, such as NeurIPS, AAAI, ACM MM, Pattern Recognition, IEEE Transactions on Cybernetics and Neurocomputing. He also serves as a reviewer for several international journals, including the IEEE Transactions on Cybernetics, the IEEE Transactions on Automation Science and Engineering, Artificial Intelligence Review, and International Journal of Machine Learning and Cybernetics.



Xiping Hu received the Ph.D. degree from the University of British Columbia, Vancouver, BC, Canada. He is currently a professor with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences and Lanzhou University, China. He was the Co-Founder and CTO of Bravolol Ltd., Hong Kong, a leading language learning mobile application company with over 100 million users, and listed as the top 2 language education platform globally. He has more than 100 papers published and presented in prestigious conferences and journals, such as IEEE TMC/TPDS/TIP/JSAC/IoT journal, IEEE COMST, IEEE COMMUNICATIONS MAGAZINE, MobiCom, AAAI, IJCAI, and WWW. He has been serving as the lead guest editors of IEEE Internet of Things Journal, IEEE Transactions on Automation Science and Engineering, and WCMC. His research areas consist of mobile cyber-physical systems, crowdsensing and affective computing.



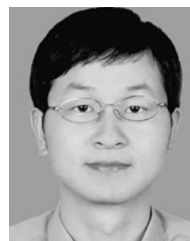
Mingkui Tan is currently a professor with the School of Software Engineering at South China University of Technology. He received his Bachelor Degree in Environmental Science and Engineering in 2006 and Master degree in Control Science and Engineering in 2009, both from Hunan University in Changsha, China. He received the Ph.D. degree in Computer Science from Nanyang Technological University, Singapore, in 2014. From 2014–2016, he worked as a Senior Research Associate on computer vision

in the School of Computer Science, University of Adelaide, Australia. His research interests include machine learning, sparse analysis, deep learning and large-scale optimization.



Yi Guo received his Ph.D. degree from the University of Greifswald, Germany, in 1997. He is currently the chief of Neurology in the Second Clinical Medical College of Jinan University, a member of the cerebrovascular disease group of the Chinese Medical Association neurology branch, and the chairman of the Shenzhen Medical Association of Neurology and the Shenzhen Medical Association of Psychosomatic Medicine. His major research areas are cerebrovascular diseases, dementia, sleep disorder, depression,

and anxiety.



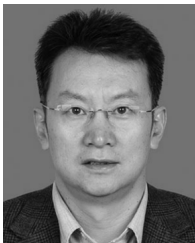
Jun Cheng received the B.Eng. and M.Eng. degrees from the University of Science and Technology of China, Hefei, China, in 1999 and 2002, respectively, and the Ph.D. degree from the Chinese University of Hong Kong, Hong Kong, in 2006. He is currently a Professor and the Founding Director of the Laboratory for Human Machine Control, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China. He has authored or coauthored about 110 articles. His current research interests include computer visions, robotics, and machine intelligence and control.

search interests include computer visions, robotics, and machine intelligence and control.



Xinwang Liu received his PhD degree from National University of Defense Technology (NUDT), China. He is now Professor of School of Computer, NUDT. His current research interests include kernel learning and unsupervised feature learning. Dr. Liu has published 60+ peer-reviewed papers, including those in highly regarded journals and conferences such as IEEE

IEEE Transactions on Pattern Analysis and Machine Intelligence, IEEE Transactions on Knowledge and Data Engineering, IEEE Transactions on Image Processing, IEEE Transactions on Neural Networks and Learning Systems, IEEE Transactions on Multimedia, IEEE Transactions on Information Forensics and Security, NeurIPS, ICCV, CVPR, AAAI, IJCAI, etc.



Bin Hu (M'05-SM'10) is currently a Professor of the School of Information Science and Engineering, Lanzhou University, Adjunct Professor with Tsinghua University, Beijing, China, and Guest Professor with ETH Zurich, Zurich, Switzerland. He is also IET Fellow, Co-Chairs of IEEE SMC TC on Cognitive Computing, and Member at Large of ACM China, Vice President of International Society for Social Neuroscience (China committee), etc. His work has been funded as a PI by the Ministry of Science and Technology,

National Science Foundation China, European Framework Programme 7, EPSRC, and HEFCE UK, etc., also published more than 100 papers in peer-reviewed journals, conferences, and book chapters including Science, Journal of Alzheimer's Disease, PLoS Computational Biology, IEEE TRANSACTION ON INTELLIGENT SYSTEMS, AAAI, BIBM, EMBS, CIKM, ACM SIGIR, etc. He has served as Chairs/Co-Chairs in many IEEE international conferences/workshops, and Associate Editors in peer-reviewed journals on Cognitive Science and Pervasive Computing, such as IEEE TRANSACTION ON AFFECTIVE COMPUTING, Brain Informatics, IET Communications, Cluster Computing, Wireless Communications, and Mobile Computing, The Journal of Internet Technology, Wiley's Security and Communication Networks, etc.