# Augmented Skeleton Based Contrastive Action Learning with Momentum LSTM for Unsupervised Action Recognition

Haocong Rao, Shihao Xu, Xiping Hu, Jun Cheng, Bin Hu

Fig. 1. Unsupervised contrastive learning paradigm for action recognition.

*Abstract*—Action recognition via 3D skeleton data is an emerging important topic. Most existing methods rely on hand-crafted descriptors to recognize actions, or perform supervised action representation learning with massive labels. In this paper, we for the first time propose a contrastive action learning paradigm named AS-CAL that exploits different augmentations of *unlabeled* skeleton sequences to learn action representations in an *unsupervised* manner. Specifically, we first propose to contrast similarity between augmented instances of the input skeleton sequence, which are transformed with multiple novel augmentation strategies, to learn inherent action patterns ("*pattern-invariance*") in different skeleton transformations. Second, to encourage learning the pattern-invariance with more consistent action representations, we propose a momentum LSTM, which is implemented as the momentum-based moving average of LSTM based query encoder, to encode long-term action dynamics of the key sequence. Third, we introduce a *queue* to store the encoded keys, which allows flexibly reusing proceeding keys to build a consistent dictionary to facilitate contrastive learning. Last, we propose a novel representation named Contrastive Action Encoding (CAE) to represent human's action effectively. Empirical evaluations show that our approach significantly outperforms hand-crafted methods by 10-50% Top-1 accuracy, and it can even achieve superior performance to many supervised learning methods[1].

*Index Terms*—Skeleton based action recognition, skeleton data augmentation, unsupervised deep learning, contrastive learning, momentum LSTM.

## I. INTRODUCTION

Human action recognition plays a vital role in computer vision. Recent development of depth sensors [1] revolutionizes the way to recognize actions, which shifts from using RGB images [2] to using depth images [3]–[7] or skeletons [8], [9]. The 3D skeleton based models [8], [10] have gained surging popularity in these years. By leveraging three-dimensional coordinates of numerous key body joints to perform action recognition, 3D skeleton based models enjoy many merits like high robustness to variations of positions, scales, and viewpoints [11].

Most existing skeleton-based methods [8], [10] utilize supervised learning paradigms to learn action representations, where an enormous number of annotations for action frames or videos are indispensable. However, labeling for a large scale dataset requires tremendous human workforce, which is usually expensive and non-scalable for many action recognition
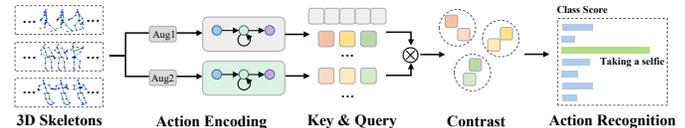
related applications [12]. In addition, there exist some challenging issues deriving from the process of manual annotation: High inter-class similarity between actions usually leads to an uncertain labeling or even mislabeling on action samples [13]. Under this circumstance, devising an effective method to learn action representations from unlabeled data attracts increasing attention [14].

Most recently, some works [3], [11], [15], [16] explore unsupervised methods to learn action features from unlabeled data. For example, [3] learns action representation via identifying correct temporal order of sequences based on an AlexNet architecture [17]. Most of these methods [11], [15], [16] rely on different paradigms of encoder-decoder [18] or generative models [19] to learn action features with a pretext task of sequential reconstruction or prediction. However, designing a pretext task to learn the data representation or distribution as losslessly as possible (*e.g.*, reconstruction) is not always sufficient for the downstream task [20]. For long action sequences with rich spatio-temporal information, it is important to keep "good" features like key action patterns and throw away trivial or noisy information to achieve a compact representation, which inherently requires an effective contrasting and learning mechanism [21].

To address all above challenges, we propose a novel **unsupervised** approach named Augmented Skeleton based Contrastive Action Learning (AS-CAL) with *momentum* Long Short-term Memory (mLSTM). AS-CAL only requires **unlabeled** 3D skeleton data to learn an effective action representation, which maximizes agreement between different augmented instances of the same action sequence with a contrastive loss. Specifically, we first propose different novel augmentation strategies to introduce specific transformations and random data perturbations to the original action sequence. Since the property of "pattern-invariance" leads to similar action patterns in a particular random transformation (*e.g.*, random rotation), we expect our model to incorporate such inherent similarity into the action representation by contrastive learning. Second, as shown in Fig. 1, given *query* and *key*

---

*(Haocong Rao and Shihao Xu contributed equally to this work.) (Corresponding authors: Xiping Hu; Jun Cheng; Bin Hu.)*

[1]Our codes are available at https://github.com/Mikexu007/AS-CAL.

sequences generated by the same skeleton data augmentation strategy or strategy composition (note that *query* and *key* sequences are randomly augmented with the same strategy, namely "Aug1" and "Aug2" shown in Fig. 1), we exploit Long Short-Term Memory (LSTM) [22] to encode query sequence, and propose a momentum LSTM (mLSTM) as the key encoder, which is implemented as momentum-based moving average of the query encoder to achieve more consistent action encoding. In this way, they encode long-term action dynamics of pairwise augmented instances (query and key) to yield the preliminary action representation for contrasting. Third, to obtain a manageable and more consistent dictionary for contrasting training samples, we introduce a *queue-based dictionary* to store keys by enqueueing the newest mini-batch keys and dequeueing the oldest ones during training, which allows our model to flexibly reuse keys from preceding mini-batches to facilitate contrastive learning. Last, we employ the contrastive loss based on Noise Contrastive Estimation (NCE) [23] to learn similarity between the query representation and positive key representation (*i.e.*, augmented instance of the same input sequence), and encourage capturing distinct action features by discriminating positive key representations from negative ones. We average action features learned from the query encoder across all time steps, and construct the final representation named Contrastive Action Encoding (CAE) for the downstream task of action recognition. We demonstrate that CAE, which is learned without any skeleton label, can be directly applied to action recognition task and achieves highly competitive performance.

In summary, we make the following contributions:

- We propose a generic unsupervised contrastive action learning paradigm named AS-CAL for action recognition. The proposed AS-CAL performs contrastive learning on action patterns of augmented skeleton sequences, which enables us to learn effective action representations from unlabeled skeleton data.
- We devise novel skeleton data augmentation strategies to generate the query and key skeleton sequences as augmented instances for contrastive learning, and showcase their effectiveness on unsupervised action representation learning.
- We propose a momentum LSTM as the key encoder with a momentum-based update of encoder's parameters, so as to enable learning more consistent action representations and facilitate contrastive action learning.
- We propose a novel action representation named Contrastive Action Encoding (CAE), which is shown to be highly effective on the action recognition task.

We comprehensively evaluate the effectiveness of our approach on four public datasets: NTU RGB+D 60 [24], NTU RGB+D 120 [25], SBU [26], and UWA3D datasets [27]. Under the linear evaluation protocol, the proposed AS-CAL significantly outperforms existing hand-crafted methods by up to $50\%$ Top-1 accuracy, and it also achieves superior performance to many supervised learning methods on NTU RGB+D 60 and NTU RGB+D 120 datasets. On SBU and UWA3D datasets, our approach is shown to perform better

than most existing supervised learning baselines.

The rest of this paper is organized as follows: Sec. II introduces relevant works and the ideas that inspire our work. Sec. III elucidates each module of the proposed approach. Sec. IV presents the details of experiments, and extensively compares our approach with existing methods. Sec. V provides ablations studies and comprehensive discussion on the proposed approach. Sec. VI draws the conclusion of this paper.

## II. RELATED WORK

In this section, we provide a comprehensive introduction for previous works in the fields of action recognition and contrastive learning.

### A. Action Recognition

Existing action recognition methods can be divided into three categories, namely (1) hand-crafted methods, (2) supervised methods, and (3) unsupervised methods. In this part, we first introduce representative hand-crafted and supervised methods in the literature (see Sec. II-A1). Then, we review state-of-the-art unsupervised methods and highlight key differences between our method and existing methods (see Sec. II-A2).

*1) Hand-Crafted and Supervised Methods:* Hand-crafted descriptors [5]–[7], [28], [29] are widely used to perform action recognition. Evangelidis *et al.* [28] design Skeletal Quads for encoding local position of joint quadruples to obtain view-invariant action features. In [5], Oreifej *et al.* use a modified histogram of oriented gradients (HOG) algorithm to extract discriminative features for action recognition. Motivated by the remarkable success achieved by recent deep neural networks (DNNs), numerous works [8]–[10] adopt DNNs to perform supervised action representation learning. By modeling the skeleton as a graph, Yan *et al.* [10] propose spatial-temporal graph convolutional networks (ST-GCN) to extract unique pattern features of different actions. Si *et al.* [9] further incorporate graph convolutional networks into LSTM to better capture discriminative features in spatial configuration and temporal dynamics for action recognition. However, these methods unexceptionally require massive labels or fine-grained annotations and cannot learn an effective action representation directly from unlabeled skeleton data.

*2) Unsupervised Methods:* Unsupervised action representation learning is a newly-emerging topic in these years. In the field of RGB-based action recognition, Srivastava *et al.* [30] propose an LSTM-based auto-encoder to learn action representations by reconstructing input videos. In [31], a hierarchical dynamic parsing and encoding method is established to model local and global temporal dynamics of action representations. Ahsan *et al.* [32] train a network to learn action features by solving the pretext task of jigsaw puzzles based on pixel patches of action sequences. Some works [4] combine depth images with RGB data to predict 3D motions and learn the view-invariant action representations. As to skeleton-based action recognition, few works like [11], [15] apply unsupervised learning to extracting unique action features from

3D skeleton data. In [11], Zheng *et al.* propose a generative adversarial network (GAN) based encoder-decoder for sequential reconstruction, and exploit the learned intermediate representation to recognize different actions. Su *et al.* [15] propose a decoder-weakening strategy for the encoder-decoder model, so as to drive the encoder to learn discriminative action features. Lin *et al.* [16] devise multiple pretext tasks (*e.g.,* motion prediction, identifying temporal order) to drive the unsupervised learning with encoder-decoder architecture, and combine them to encourage the Bi-GRU encoder to capture more action patterns.

There are a few key differences between our method and previous skeleton-based methods: (1) We propose a novel contrastive action learning paradigm to learn effective action representations from unlabeled skeleton data. We do NOT require feature engineering like [5], [6], [28] or designing task-specific models (*e.g.*, GAN [11], encoder-decoder [15]) to implement corresponding pretext tasks like reconstruction. By contrast, we exploit different novel skeleton data augmentation strategies to drive the contrastive learning, which encourages the model to learn inherent action patterns from different skeleton transformations. Besides, the proposed contrastive learning paradigm is highly flexible and scalable, which could be extended to different pretext tasks and encoders. (2) The property of consistence is exploited to achieve better contrastive learning: We not only propose a momentum LSTM to learn more consistent action representations but also involve a queue to build a consistent and memory-efficient dictionary to improves the performance of the unsupervised contrastive action learning.

### B. Contrastive Learning

Contrastive learning [20], [23], [33]–[36] is an effective unsupervised learning method that can be applied on various pretext tasks via a contrastive loss. Pretext tasks (*e.g.*, motion reconstruction, frame prediction) can be used to learn useful data representation beforehand and later to be applied to the tasks of real interest like action recognition. Some works design pretext tasks based on auto-encoders to denoise images [37], or achieve plausible image colorization [38]. The contrastive loss [39] is associated with tasks and it measures the similarities of sample pairs in a representation space. For example, in instance discrimination task [33], the noise-contrastive estimation (NCE) related contrastive loss [40] pulls closer the augmented samples from the same instance, and pushes apart ones from different instances. The contrastive multiview coding (CMC) [20] aims to maximize mutual information between different views, while the momentum contrastive paradigm (MoCo) [34], [36] facilitates contrastive unsupervised learning by queue-based dictionary look-up mechanism and the momentum-based update. Compared with Moco, SwAV [41] incorporates the online clustering into contrastive learning, which runs with small batches and requires less memory for storage of features. In [35], Chen *et al.* propose SimCLR with the multi-layer perceptron (MLP) projection head and stronger color augmentation to further improve the quality of unsupervised learned representation (Likewise, Moco v2 [36]
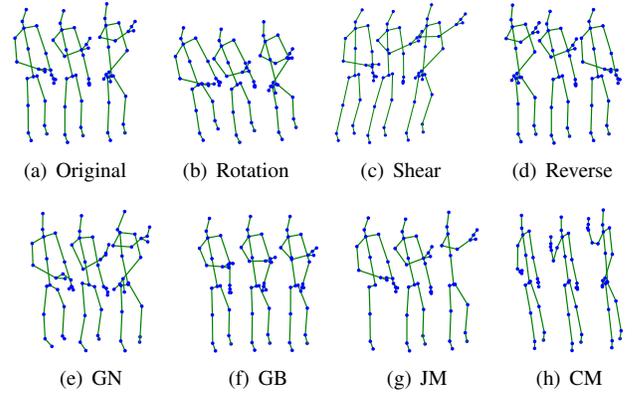


(a) Original    (b) Rotation    (c) Shear    (d) Reverse

(e) GN    (f) GB    (g) JM    (h) CM

Fig. 2. Visualization of data augmentations (b)-(h) for the same skeleton sequence (a).

benefits from the MLP and the stronger color augmentation). SimCLRv2 [42] rather adopts larger ResNet models, deeper projection heads, and memory mechanism from MoCo to achieve superior performance. In a simpler way, SimSiam [43] requires neither negative pairs nor a momentum encoder (*i.e.*, can be viewed as "SimCLR" without negative pairs) to obtain competitive outcomes. [34]–[36], [42], [43] could be viewed as an instance discrimination method to perform unsupervised visual representation learning. This work is the first attempt to explore contrastive learning based on instance discrimination for learning an effective action representation directly from unlabeled 3D skeleton sequences.

### III. PROPOSED APPROACH

Suppose that an input skeleton sequence $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_T)$ contains $T$ consecutive skeleton frames, where $\boldsymbol{x}_i \in \mathbb{R}^{M \times J \times 3}$ contains 3D coordinates of $J$ different body joints for $M$ actors. The training set $\Phi = \{\boldsymbol{x}^{(i)}\}_{i=1}^{N}$ contains $N$ skeleton sequences of different actions collected from multiple views and persons. Each skeleton sequence $\boldsymbol{x}^{(i)}$ corresponds to a label $y_i$, where $y_i \in \{a_1, \cdots, a_c\}$, $a_i$ represents the $i^{th}$ action class, and $c$ is the number of action classes. Our goal is to learn an effective action representation $\boldsymbol{q}$ from $\boldsymbol{x}^{(i)}$ without using any skeleton label. Then, the effectiveness of learn features $\boldsymbol{q}$ is validated by the linear evaluation protocol: Leaned features $\boldsymbol{q}$ and labels are used to train a linear classifier for action recognition (note that $\boldsymbol{q}$ is frozen and NOT tuned at the recognition stage). The overview of the proposed approach is given in Fig. 3, and we present the details of each technical component below.

### A. Data Augmentation for Skeleton Sequences

As the goal of contrastive learning is to learn shared pattern information between different augmented instances of the same example [34]–[36], it is natural to consider the property of "pattern-invariance" in skeleton sequences: Random transformations of the same skeleton sequence under a specific augmentation strategy or strategy composition (*e.g.*, rotation) KEEP similar action patterns, which can be contrasted and learned to achieve an effective action representation. Take

action selfie as an example, when we randomly rotate the corresponding skeleton sequence in a specific direction, we could still easily recognize this action since the similar hand pose remains regardless of different rotation angles. To better learn such "pattern-invariance" and yield a robust action representation, we introduce appropriate data perturbations to randomly transform the skeleton sequence, and perform contrastive learning to encode the inherently similar action pattern. As presented by Fig. 2, we devise seven augmentation strategies containing 3D transformations of skeleton sequences, with definitions as below:

(1) *Rotation* (see Fig. 2 (b)). The Euler's rotation theorem ensures that any 3D rotation can be composed of rotations about three axes [44]. The three basic rotation matrices with rotate angles $\alpha, \beta, \gamma$ on $X, Y, Z$ axis respectively are given as follows:

$$\boldsymbol{R}_X(\alpha) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\alpha & -\sin\alpha \\ 0 & \sin\alpha & \cos\alpha \end{bmatrix} \quad (1)$$

$$\boldsymbol{R}_Y(\beta) = \begin{bmatrix} \cos\beta & 0 & \sin\beta \\ 0 & 1 & 0 \\ -\sin\beta & 0 & \cos\beta \end{bmatrix} \quad (2)$$

$$\boldsymbol{R}_Z(\gamma) = \begin{bmatrix} \cos\gamma & -\sin\gamma & 0 \\ \sin\gamma & \cos\gamma & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (3)$$

$$\boldsymbol{R} = \boldsymbol{R}_Z(\gamma)\boldsymbol{R}_Y(\beta)\boldsymbol{R}_X(\alpha) \quad (4)$$

where $\boldsymbol{R}_X(\alpha), \boldsymbol{R}_Y(\beta), \boldsymbol{R}_Z(\gamma)$ indicate the rotation matrices on $X, Y, Z$ axis with random angles $\alpha, \beta, \gamma$ respectively, and $\boldsymbol{R}$ is a general rotation matrix obtained from three basic rotation matrices (Eqn. 1, 2, 3) using matrix multiplication.

To simulate the viewpoint changes of the camera on each axis, we design the following rotation strategy: For all joint coordinates in a skeleton sequence, we randomly choose a main rotation axis $A \in \{X, Y, Z\}$ and select a random rotation angle from $[0, \frac{\pi}{6}]$ for the axis $A$, while the remaining two axes perform rotations with random angles from $[0, \frac{\pi}{180}]$, which aims to introduce random rotation perturbations to improve the robustness of our model to viewpoint changes [44]. Then, we apply the rotation $\boldsymbol{R}$ to original coordinates of the skeleton sequence and get the transformed coordinates.

(2) *Shear* (see Fig. 2 (c)). The shear transformation is a linear mapping matrix that displaces each joint in a fixed direction, *i.e.*, the shape of 3D coordinates of body joints will be slanted with a random angle. We define the shear transformation matrix as follows:

$$\boldsymbol{S} = \begin{bmatrix} 1 & s_X^Y & s_X^Z \\ s_Y^X & 1 & s_Y^Z \\ s_Z^X & s_Z^Y & 1 \end{bmatrix} \quad (5)$$

where $s_X^Y, s_X^Z, s_Y^X, s_Y^Z, s_Z^X, s_Z^Y \in [-1, 1]$ are randomly sampled shear factors from each dimension to another (*e.g.*, $s_X^Y$ corresponds to the shear factor from $X$ to $Y$). We transform all joint coordinates of the original skeleton sequence with the shear matrix.

(3) *Reverse* (see Fig. 2 (d)). Similar to the operation of horizontal flip in image augmentation, we consider "flip" from

the view of temporal order: The order of original skeleton sequence is reversed at 50% chance. Inspired by the fact that the order of skeleton sequence may NOT influence human's perception of actions, we expect the model to learn crucial action details (*e.g.*, joint positions, joint angles) from a reverse sequence.

(4) *Gaussian Noise (GN)* (see Fig. 2 (e)). To simulate the noisy positions caused by estimation or annotation, we add Gaussian noise $\mathcal{N}(0, 0.05)$ over joint coordinates of the original sequence.

(5) *Gaussian Blur (GB)* (see Fig. 2 (f)). As an effective augmentation strategy to reduce the level of details and noise of images, Gaussian blur can be applied to the skeleton sequence to smooth noisy joints and decrease action details. We randomly sample $\sigma \in [0.1, 2.0]$ for the Gaussian kernel, which is a sliding window with length of 15. Joint coordinates of the original sequence are blurred at 50% chance by the kernel $G(\cdot)$ below:

$$G(t) = \exp(-\frac{t^2}{2\sigma^2}), \quad t \in \{-7, -6, \cdots, 6, 7\}, \quad (6)$$

where $t$ denotes the relative position from the center skeleton (note that the negative/positive number indicates the skeletons before/after the center skeleton), and the length of the kernel is set to 15 corresponding to the total span of $t$.

(6) *Joint Mask (JM)* (see Fig. 2 (g)). We employ a zero-mask to a number of body joints in skeleton frames (*i.e.*, replace all coordinates by zeros), which encourages the model to learn different local regions (*i.e.*, except for the masked region) that probably contain crucial action patterns. To be more specific, we randomly choose a certain number of body joints (number of joints $\overline{V} \in \{5, 6, \cdots, 15\}$) from random frames (number of frames $\overline{L} \in \{50, 51, \cdots, 100\}$) in the original skeleton sequence to apply the zero-mask.

(7) *Channel Mask (CM)* (see Fig. 2 (h)). We randomly choose a "channel" (*i.e.*, an axis $A \in \{X, Y, Z\}$) of skeleton sequence, and apply a zero mask to all coordinates on this axis. In this way, the original skeleton sequence can be transformed to 2D projection sequence, which enables the model to learn dominant action patterns from a particular plane.

**Remark:** To sample the query and key for the same sequence $\boldsymbol{x}$, we adopt the same augmentation strategy or strategy composition to randomly transform $\boldsymbol{x}$ to query sequence $\tilde{\boldsymbol{x}}$ and key sequence $\overline{\boldsymbol{x}}$, which are then fed into the encoders for action dynamics learning. We demonstrate that the proposed augmentation strategies improve the performance of both contrastive learning and action recognition (see Sec. V-B).

### B. Augmented Skeleton based Contrastive Action Learning (AS-CAL)

The nature of "pattern-invariance" endows randomly augmented instances of the same skeleton sequence with highly similar patterns, which allows the model to learn good representations by contrasting the similarity between sequence's different transformations. To this end, we propose an unsupervised learning approach named Augmented Skeleton based Contrastive Action Learning (AS-CAL) with a *momentum* LSTM and a *queue-based dictionary* to maximize inherent
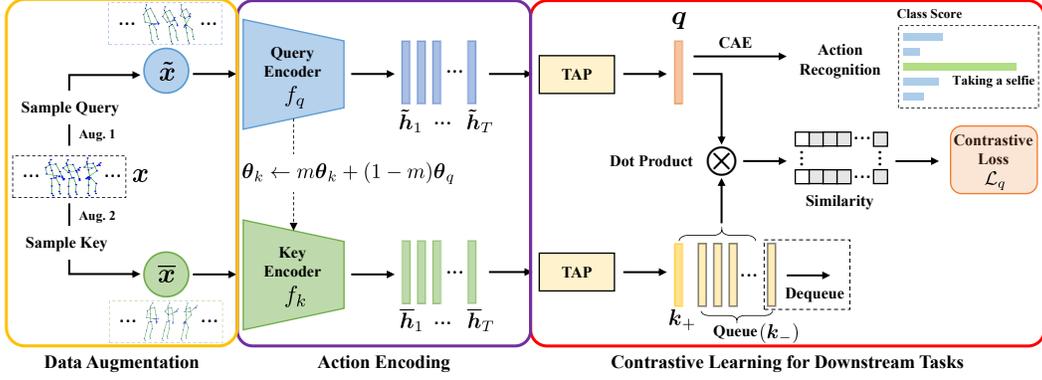
Fig. 3. Flow diagram of the proposed AS-CAL (detailed in Sec. III-B1).

agreement between different augmented instances of the same skeleton sequence via an effective contrastive loss, so as to learn an effective action representation for the action recognition task. In this section, we first provide an overall description for the working flow of proposed AS-CAL (see Sec. III-B1). Then, we systematically illustrate its technical components (see Sec. III-B2, Sec. III-B3, Sec. III-B4, and Sec. III-B5). In Sec. III-C and Sec. III-D, we elaborate the proposed Contrastive Action Encodings (CAE) and summarize our approach in the form of an algorithm (see 1 and computation flow.

*1) Working Flow of AS-CAL:* The overview of proposed AS-CAL is shown in Fig. 3, and its working flow can be described by five steps: **(1)** First, we sample query $\tilde{x}$ and key $\overline{x}$ from the input skeleton sequence $x$ by two random augmentations using the same strategy or strategy composition (see yellow box in Fig. 3). **(2)** Second, the query encoder $f_q$ and the momentum-based key encoder $f_k$, which updates its parameters $\boldsymbol{\theta}_k$ with weighted average of $m\boldsymbol{\theta}_k$ and $(1-m)\boldsymbol{\theta}_q$, encode skeleton frames of $\tilde{x}$ and $\overline{x}$ into hidden states $\tilde{h}$ and $\overline{h}$ to represent action encoding information (see purple box in Fig. 3). **(3)** Third, all hidden states are then averaged across time (TAP) to obtain query representation $\boldsymbol{q}$ and positive key representation $\boldsymbol{k}_+$. **(4)** Then, the oldest batch of negative keys ($\boldsymbol{k}_-$) in queue is dequeued while the new batch of $\boldsymbol{k}_+$ is enqueued. **(5)** Finally, dot products between query and all keys are computed, and the similarity between positive pairs is maximized by contrastive loss $\mathcal{L}_q$ (see red box in Fig. 3). The learned Contrastive Action Encoding (CAE) $\boldsymbol{q}$ is fed into a linear classifier for action recognition. We present the motivation and technical details for each technical component of AS-CAL below.

*2) Momentum LSTM (mLSTM):* Larger dictionary providing more negative keys can help achieve better contrastive learning with higher training efficiency (*e.g.*, faster convergence [35]). However, it is usually intractable for the key encoder to update its parameters with all samples in the large dictionary [34]. To perform long-term action dynamics learning and keep key representations' consistency better, we exploit an LSTM (denoted as $f_q$) to encode the query sequence $\tilde{x}$, and propose a momentum LSTM (mLSTM) as the key encoder (denoted as $f_k$): The mLSTM does NOT perform back-propagation but updates its parameter $\boldsymbol{\theta}_k$ by
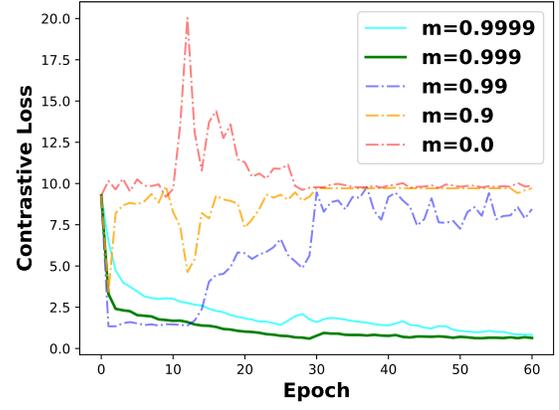


Fig. 4. Contrastive loss curves of the proposed AS-CAL with different momentum coefficients $m$ on the NTU RGB+D 60 dataset (C-Sub).

the exponentially weighted average (*i.e.*, momentum-based moving average) of its original parameters $\boldsymbol{\theta}_k$ and parameters $\boldsymbol{\theta}_q$ of the query encoder.

Formally, given a query sequence $\tilde{x}$ and a key sequence $\overline{x}$ generated by augmentations of the input skeleton sequence, we use the query encoder $f_q$ and key encoder $f_k$, which are build with the LSTM, to encode each skeleton frame into hidden states as follows (see Fig. 3):

$$\tilde{h}_t = \begin{cases} f_q\left(\tilde{x}_1\right) & \text{if } t = 1 \\ f_q\left(\tilde{h}_{t-1}, \tilde{x}_t\right) & \text{if } t > 1 \end{cases} \qquad (7)$$

$$\overline{h}_t = \begin{cases} f_k\left(\overline{x}_1\right) & \text{if } t = 1 \\ f_k\left(\overline{h}_{t-1}, \overline{x}_t\right) & \text{if } t > 1 \end{cases} \qquad (8)$$

where $\tilde{h}_t, \overline{h}_t \in \mathbb{R}^E$, $t \in \{1, \cdots, T\}$ denotes the skeleton frame number. $\tilde{h}_1, \cdots, \tilde{h}_T$ and $\overline{h}_1, \cdots, \overline{h}_T$ are encoded hidden states of the query sequence $\tilde{x}_1, \cdots, \tilde{x}_t$ and key sequence $\overline{x}_1, \cdots, \overline{x}_t$ respectively, and they contain preliminary action encoding information. In the **training** stage, the key encoder $f_k$ is an mLSTM with its parameters updated as below:

$$\boldsymbol{\theta}_k \leftarrow m\boldsymbol{\theta}_k + (1-m)\boldsymbol{\theta}_q \qquad (9)$$

where $\boldsymbol{\theta}_k$, $\boldsymbol{\theta}_q$ are parameters of the key encoder $f_k$ and query encoder $f_q$ respectively, $m \in [0,1)$ is a momentum coefficient to control the update speed. The momentum-based update

makes $\theta_k$ evolves more smoothly than $\boldsymbol{\theta}_q$ (note that ONLY the parameters of query encoder ($\boldsymbol{\theta}_q$) are updated by back-propagation while the parameters of key encoder ($\boldsymbol{\theta}_k$) are updated by Eqn. 9). In this way, the difference among key encoders remains very small at different iterations (*i.e.*, in different mini-batches), which encourages the model to keep consistency of key representations ($\overline{\boldsymbol{h}}$). Compared with the full or fast update ($m \rightarrow 0$) that undergoes drastic parameters' changes, the lower evolving speed ($0.999 \leq m < 1$) of key encoder benefits contrastive learning, which is demonstrated by the Fig. 4: When $m \leq 0.99$, contrastive loss curves show more fluctuations as $m$ gets smaller, and the model fails to converge to a low loss stably. In contrast, the proposed AS-CAL with $m > 0.99$ can achieve an evidently lower contrastive loss with a faster convergence, and $m = 0.999$ is shown to be the best performer. In Sec. V-A3, we demonstrate that a better contrastive learning (AS-CAL) encourages learning a more effective action representation for action recognition.

*3) Temporal Average Pooling:* Temporal average pooling (TAP) is the implementation of average pooling in the temporal domain, which can be used to aggregate global action encoding information across time [9]. In this work, we apply TAP to hidden states of $\tilde{\boldsymbol{x}}$ and $\overline{\boldsymbol{x}}$ to yield the query representation $\boldsymbol{q}$ and corresponding positive key representation $\boldsymbol{k}_+$ below:

$$\boldsymbol{q} = \text{TAP}(\tilde{\boldsymbol{h}}_1, \cdots, \tilde{\boldsymbol{h}}_T) = \frac{1}{T}\sum_{i=1}^{T} \tilde{\boldsymbol{h}}_i \qquad (10)$$

$$\boldsymbol{k}_+ = \text{TAP}(\overline{\boldsymbol{h}}_1, \cdots, \overline{\boldsymbol{h}}_T) = \frac{1}{T}\sum_{i=1}^{T} \overline{\boldsymbol{h}}_i \qquad (11)$$

where $\boldsymbol{q}$, $\boldsymbol{k}_+ \in \mathbb{R}^E$, $\tilde{\boldsymbol{h}}_i$ and $\overline{\boldsymbol{h}}_i$ are the $i^{th}$ hidden states of $\tilde{\boldsymbol{x}}$ and $\overline{\boldsymbol{x}}$ respectively. $\boldsymbol{k}_+$ represents the positive key representation corresponding to the query representation $\boldsymbol{q}$. As shown in Fig. 3, at each training step, a new batch of pairwise query and key sequences are encoded and pooled into $\boldsymbol{q}$ and $\boldsymbol{k}_+$ for contrastive learning.

*4) Queue-Based Dictionary:* To build a large, consistent, and manageable dictionary for AS-CAL, we introduce a queue of size $K$ to maintain encoded keys: At each training step, the current mini-batch of keys is enqueued to the dictionary while the oldest mini-batch of keys in the queue is removed (see Fig. 3). This allows us to progressively replace the samples in the dictionary and reuse the preceding encoded keys (note that all preceding keys in the queue are viewed as negative keys ($\boldsymbol{k}_-$) in the training of new mini-batch). As presented in Table I, we compare three generic contrastive learning paradigms: (a) The proposed AS-CAL using queue-based dictionary. (b) The end to end paradigm using mini-batch based dictionary without momentum-based encoder [21], [39], [45]. (c) The memory bank paradigm [33] with momentum update on representations of the same sample (with no key encoder). Compared with other dictionary structures (mini-batch or memory bank based dictionary), the queue-based dictionary has several prominent advantages: (1) Using the queue can build a flexible and much larger dictionary than a typical mini-batch, whose dictionary size is limited by the device memory and the large-batch optimization [46]. (2) The

| Paradigm | Dictionary Size | Back-Propagation | Sampling Source |
|---|---|---|---|
| AS-CAL (Queue) | Size of queue ($K$) | Only $f_q$ requires | Current queue |
| End to End | Size of mini-batch ( Typically $< K$) | Both $f_q$ and $f_k$ require | Current batch |
| Memory Bank | Size of all samples (Typically $\gg K$) | No $f_k$, only $f_q$ requires | Memory bank of past epoch |

queue is more memory-efficient than the memory bank that stores all keys of the dataset. Meanwhile, the memory bank only samples keys from the past epoch, while the queue maintains the immediate mini-batches of keys to achieve more consistent dictionary. Quantitative results and analysis are in Sec. V-C, and we demonstrate that the proposed AS-CAL with queue-based dictionary can achieve superior performance to existing contrastive paradigms.

*5) Contrastive Loss:* As the goal of AS-CAL is to learn an effective representation of inherent action patterns by contrasting different transformations of skeleton sequences, we expect the model to maximize the similarity between augmented instances of the same sequence: $\boldsymbol{q}$ and its matched positive key $\boldsymbol{k}_+$ are supposed to be similar while the dissimilar ones ($\boldsymbol{q}$ and negative keys in queue) should be separated. We use dot product to measure the similarity and employ the contrastive loss function InfoNCE [23] to perform AS-CAL:

$$\mathcal{L}_q = -\log \frac{\exp\left(\boldsymbol{q} \cdot \boldsymbol{k}_+/\tau\right)}{\exp\left(\boldsymbol{q} \cdot \boldsymbol{k}_+/\tau\right) + \sum_{i=1}^{K} \exp\left(\boldsymbol{q} \cdot \boldsymbol{k}_-^i/\tau\right)} \qquad (12)$$

where $\mathcal{L}_q$ denotes the contrastive loss, $\tau$ is a temperature hyper-parameter to adjust the contrastive learning, $K$ is the number of keys in the queue, and $\boldsymbol{k}_-^i$ is the $i^{th}$ negative key in the queue. The main algorithm of AS-CAL is presented in Algorithm 1.

### C. Contrastive Action Encoding (CAE)

Since our ultimate goal is to learn good action features from skeleton data to perform action recognition, we need to extract certain internal embedding of skeleton sequences from the proposed AS-CAL as the final action representation. We recall that the query encoder $f_q$ drives the momentum update of the key encoder $f_k$, and it can encode the long-term action dynamics of the skeleton sequence to achieve an effective action representation for contrastive learning. Hence, we use the $f_q$ learned by the proposed AS-CAL as the final action encoder. The pre-trained $f_q$ (note that $f_q$ is frozen in the linear evaluation stage) encodes the original skeleton sequence $\boldsymbol{x}$ into hidden states, and applies TAP to yield the final action representation named Contrastive Action Encoding (CAE) for action recognition. By combining Eqn. 7 and Eqn. 10, we give the computation of CAE as follows:

$$\boldsymbol{q} = \text{TAP}(f_q(\boldsymbol{x})) \qquad (13)$$

where $f_q$ is pre-trained by the proposed AS-CAL, and $\boldsymbol{q}$ is the CAE that aggregates the global action encoding information in an average manner. Here we use the same symbol $\boldsymbol{q}$ (*i.e.*, same as the query action representation of transformed sequence in Eqn. 10) to represent CAE because $\boldsymbol{x}$ can be viewed as

---

**Algorithm 1** Main algorithm of AS-CAL

---

**Input:** Temperature $\tau$, momentum coefficient $m$, mini-batch size $n$, query encoder $f_q$, key encoder $f_k$, queue size $K$

\# Initialization

Randomly initialize parameters $\boldsymbol{\theta}_q$ of $f_q$, and copy to $f_k$ (parameters $\boldsymbol{\theta}_k$)

Randomly initialize negative keys $\left\{\boldsymbol{k}_-^j\right\}_{j=1}^K$ in queue \# $\boldsymbol{k}_-^j \in \mathbb{R}^E$

**for** a sampled mini-batch $\left\{\boldsymbol{x}^{(i)}\right\}_{i=1}^n$ **do**

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

        \# Select one or a composition of augmentation strategies to perform two random augmentations: $\text{Aug1}(\cdot)$, $\text{Aug2}(\cdot)$

        \# The first augmentation to get queries

        $\tilde{\boldsymbol{x}}^{(i)} = \text{Aug1}(\boldsymbol{x}^{(i)})$

        $(\tilde{\boldsymbol{h}}_1, \cdots, \tilde{\boldsymbol{h}}_T) = f_q(\tilde{\boldsymbol{x}}^{(i)})$          \# $\tilde{\boldsymbol{h}}_i \in \mathbb{R}^E$

        $\boldsymbol{q}^{(i)} = \text{TAP}(\tilde{\boldsymbol{h}}_1, \cdots, \tilde{\boldsymbol{h}}_T)$          \# $\boldsymbol{q}^i \in \mathbb{R}^E$

        \# The second augmentation to get positive keys

        $\overline{\boldsymbol{x}}^{(i)} = \text{Aug2}(\boldsymbol{x}^{(i)})$

        $(\overline{\boldsymbol{h}}_1, \cdots, \overline{\boldsymbol{h}}_T) = f_k(\overline{\boldsymbol{x}}^{(i)})$        \# $\overline{\boldsymbol{h}}_i \in \mathbb{R}^E$

        $\boldsymbol{k}_+^{(i)} = \text{TAP}(\overline{\boldsymbol{h}}_1, \cdots, \overline{\boldsymbol{h}}_T)$        \# $\boldsymbol{k}_+^{(i)} \in \mathbb{R}^E$

        detach $\boldsymbol{k}_+^{(i)}$          \# No gradient to keys

    **end for**

    \# Calculate contrastive loss $\mathcal{L}_q$ for mini-batch and update encoders

    $\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n \log \frac{\exp\left(\boldsymbol{q}^{(i)} \cdot \boldsymbol{k}_+^{(i)} / \tau\right)}{\exp\left(\boldsymbol{q}^{(i)} \cdot \boldsymbol{k}_+^{(i)} / \tau\right) + \sum_{j=1}^K \exp\left(\boldsymbol{q}^{(i)} \cdot \boldsymbol{k}_-^{(j)} / \tau\right)}$

    Update $f_q$ to minimize $\mathcal{L}$

    Update $f_k$ with momentum: $\boldsymbol{\theta}_k \leftarrow m\boldsymbol{\theta}_k + (1-m)\boldsymbol{\theta}_q$

    \# Update queue

    Enqueue keys of current mini-batch $\left\{\boldsymbol{k}_+^{(i)}\right\}_{i=1}^n$

    Dequeue the oldest mini-batch of keys

**end for**

---

an identity transformation of the input skeleton sequence to generate the query representation for action recognition.

**Comparison with Other Action Representations.** In this work, we explore potential action representations and evaluate their performance on the action recognition task: (1) $\tilde{\boldsymbol{h}}_T$: The final hidden state from $f_q$. (2) $\overline{\boldsymbol{h}}_T$: The final hidden state from $f_k$. (3) $\boldsymbol{k}$: The key representation from $f_k$. (4) CAE: The query representation $\boldsymbol{q}$ from $f_q$. (5) CAE+: The combination (concatenation) of $\boldsymbol{q}$ (CAE) and $\boldsymbol{k}$. We follow the linear evaluation protocol (see Sec. IV-B) to validate their effectiveness. The quantitative results are reported in the supplementary material, which demonstrates that CAE (comparable to (5)) is the best performer over other action representations (1) (2) (3).

### D. The Entire Approach

To summarize intuitively, we represent the operation and computation flow of the query representation $\boldsymbol{q}$ and positive key representation $\boldsymbol{k}_+$ as follows (see Fig. 3):

- $\boldsymbol{x} \rightarrow \text{Aug. } 1(\boldsymbol{x}) \rightarrow \tilde{\boldsymbol{x}} \rightarrow f_q(\tilde{\boldsymbol{x}}) \rightarrow \tilde{\boldsymbol{h}} \rightarrow \text{TAP} \rightarrow \boldsymbol{q}$
- $\boldsymbol{x} \rightarrow \text{Aug. } 2(\boldsymbol{x}) \rightarrow \overline{\boldsymbol{x}} \rightarrow f_k(\overline{\boldsymbol{x}}) \rightarrow \overline{\boldsymbol{h}} \rightarrow \text{TAP} \rightarrow \boldsymbol{k}_+$

Here "Aug. 1" and "Aug. 2" are two random augmentations based on the same augmentation strategy or augmentation strategy composition to transform the input skeleton sequence (note that we extensively evaluate different compositions of augmentation strategies in Sec. V-B). "$\tilde{\boldsymbol{h}}$" (Eqn. 7) and "$\overline{\boldsymbol{h}}$" (Eqn. 8) represent hidden states of the query sequence and key sequence. $\boldsymbol{q}$ (Eqn. 10) and $\boldsymbol{k}_+$ (Eqn. 11) denote the positive pair of key and query representations for $\boldsymbol{x}$. During

the training process of AS-CAL, $f_k$ applies the momentum update of parameters (Eqn. 9) following $f_q$, and the InfoNCE loss function $\mathcal{L}_q$ (Eqn. 12) guides the whole contrastive learning (see Algorithm 1). For the downstream task of action recognition, we employ the cross-entropy loss to train the linear classifier on the proposed CAE ($\boldsymbol{q}$) (Eqn. 13).

## IV. EXPERIMENTS

In this section, we conduct extensive experiments to demonstrate the effectiveness of our method on action recognition with a comprehensive comparison to the existing state-of-the-art and mainstream methods. The evaluation is performed on four public skeleton-based action datasets: (1) Two large-scale datasets: NTU RGB+D 60 Action Dataset [24] and NTU RGB+D 120 Action Dataset [25]; (2) Two small-scale datasets: SBU Kinect Interaction Dataset [26] and UWA3D Multiview Activity II [27].

### A. Dataset

**NTU RGB+D 60 Action Dataset** [24]: A large-scale action recognition dataset including 60 classes of actions collected from 40 different subjects. The total number of action skeleton sequences is 56578. Two evaluation protocols are provided: (1) Cross-Subject setting (C-Sub) that splits the dataset into training set with 40091 samples and testing set with 16487 samples. (2) Cross-View setting (C-View) that utilizes 18932 sequences recorded by one camera for testing and other 37646 ones for training.

**NTU RGB+D 120 Action Dataset** [25]: As an extension of NTU RGB+D 60 dataset, this dataset contains 120 actions from 106 subjects and the total number of action skeleton sequences is 113945. Likewise, two evaluation protocols are provided: Cross-Subject (C-Sub) and Cross-Setup (C-Set). In C-Sub setting, 63026 sequences of 53 subjects are used for training while other 50919 ones are exploited for testing. In C-Set setting, different setups are used for training (16 setups) and testing (16 setups).

**SBU Kinect Interaction Dataset (SBU)** [26]: The SBU dataset is a two-person based interaction action dataset. It contains 8 types of interactions in 282 short videos with depth images, RGB images, and 3D skeletons. Each skeleton in this dataset contains 3D coordinates of only 15 joints, and we use all skeleton sequences to train our model. We adopt the 5-fold cross-validation [26] and report the average results.

**UWA3D Multiview Activity II (UWA3D)** [27]: It consists of 30 different actions performed by 10 subjects. It provides actions samples from 4 different views: front (V1), left side (V2), right side (V3), and top view (V4). The total number of action sequences is 1075. The High inter-class similarity of actions (*e.g.,* drinking and making phone call) and diversity of views points make the action recognition very challenging.

### B. Implementation Details

Our experiments are implemented by two parts: (1) Unsupervised pre-training without using skeleton labels to learn action representation; (2) Linear evaluation on the learned

representation to validate their effective on the action recognition task. We detail these two parts along with their default configurations below.

*1) Unsupervised Pre-training:* The proposed AS-CAL approach, including the LSTM-based query encoder $f_q$ and the mLSTM-based key encoder $f_k$, is pre-trained to learn an effective action representation from *unlabeled* skeleton sequences. We opt for SGD as the optimizer with weight decay of $1e^{-4}$ and SGD momentum of 0.9. The pre-training runs for 60 epochs with an initial learning rate of 0.01, which is multiplied by 0.1 at 30 epochs.

*2) Linear Evaluation Protocol:* To validate the effectiveness of the proposed action representation CAE, we follow the linear evaluation protocol [11], [34], [35], which trains a linear classifier attached to the frozen model (note that all encoders keep the parameters learned by AS-CAL and are NOT tuned at this training stage). After training the linear classifier using skeleton sequences and corresponding labels in the training set, the effectiveness of action representations can be evaluated by the recognition accuracy on the testing set. During the linear evaluation, SGD optimizer is used with a Nesterov momentum of 0.9 and an initiate learning rate of 1. Within 90 training epochs, the learning rate is decayed by $0.5\times$ at 15, 35, 60, and 75 epochs. We report Top-1 accuracy for the linear evaluation.

*3) Default Configurations:* The sequence length $T$ is set to 150, 40, 60 for NTU RGB+D 60/120, SBU, UWA3D dataset[2] respectively. As for the actors in the sample, we select first two actors[3] ($M = 2$). We subtract the coordinate of the middle spine joint from coordinates of all joints to make a normalization of skeleton sequences. For data augmentation, we sequentially apply random reverse and shear to skeleton sequences as the default setting (illustrated in Sec. III-A). Note that data augmentation strategies are only used in the unsupervised training stage. We use two-layer LSTM with $E = 256$ hidden units per layer to build the key encoder and the query encoder on NTU RGB+D 60/120 datasets, while we use single layer LSTM with $E = 256$ hidden units on the rest of datasets (note that the representation before all projection heads is a 256-dimensional vector in Sec. V-A1). For SBU and UWA3D datasets, we implement three supervised baseline methods (one-layer RNN/GRU/LSTM with 256 hidden units) for comparison. The momentum coefficient $m$ is set to 0.999. The queue size $K$ is set to 16384, 200, and 500 for NTU RGB+D 60/120, SBU, and UWA3D datasets respectively. The temperature $\tau$ is set to 0.06. The size of mini-batch is set to 32 for all experiments.

### C. Performance Comparison

In Table II and Table III, we conduct an extensive comparison with existing supervised and unsupervised methods on two large-scale datasets (NTU RGB+D 60 and NTU RGB+D 120), and also include hand-crafted methods as a reference. For SBU (Table V) and UWA3D datasets (Table

[2]If a sample sequence is not long enough, we make zero padding.
[3]We use the default actor order of the dataset. If the actor number $M < 2$, we make zero padding.

TABLE II
COMPARISON WITH HAND-CRAFTED, SUPERVISED, AND UNSUPERVISED METHODS ON NTU RGB+D 60 DATASET. "*" REPRESENTS DEPTH IMAGE BASED METHODS. BOLD NUMBERS REFER TO THE BEST PERFORMERS.

| Id | Method | C-View Accuracy (%) | C-Sub Accuracy (%) |
|----|--------|---------------------|--------------------|
| | **Hand-Crafted Methods** | | |
| 1 | *HON4D [5] | 7.3 | 30.6 |
| 2 | *Super Normal Vector [6] | 13.6 | 31.8 |
| 3 | *HOG$^2$ [7] | 22.3 | 32.2 |
| 4 | Skeletal Quads [28] | 41.4 | 38.6 |
| 5 | Lie Group [29] | 52.8 | 50.1 |
| | **Supervised Methods** | | |
| 6 | HBRNN [47] | 64.0 | 59.1 |
| 7 | Deep RNN [24] | 64.1 | 56.3 |
| | **Unsupervised Methods** | | |
| 8 | *Shuffle&Learn [3] | 40.9 | 46.2 |
| 9 | *Li *et al..* [4] | 53.9 | **60.8** |
| 10 | LongT GAN [11] | 48.1 | 39.1 |
| 11 | MS$^2$L [16] | - | 52.6 |
| 12 | Ours (CAE) | 63.6 | 58.0 |
| 13 | Ours (CAE+) | **64.8** | 58.5 |

TABLE III
COMPARISON WITH SUPERVISED LEARNING METHODS ON NTU RGB+D 120 DATASET.

| Id | Method | C-Set Accuracy (%) | C-Sub Accuracy (%) |
|----|--------|---------------------|--------------------|
| | **Supervised Methods** | | |
| 1 | Soft RNN [48] | 44.9 | 36.3 |
| 2 | Part-Aware LSTM [24] | 26.3 | 25.5 |
| | **Unsupervised Methods** | | |
| 3 | Ours (CAE) | 49.2 | 48.3 |
| 4 | Ours (CAE+) | **49.2** | **48.6** |

IV), we compare AS-CAL with three supervised learning baselines (RNN, GRU, LSTM). In Table II, III, V, IV, we simultaneously compare the performance of the proposed CAE and its enhanced representation CAE+ (see Sec. 3.3), which is the concatenation of CAE (*i.e.*, the query representation $\boldsymbol{q}$) and the key representation ($\boldsymbol{k}$). The crucial results are reported as below.

*1) Comparison with Unsupervised Methods:* As shown in Table II, our approach enjoys distinct advantages over existing unsupervised methods (Id = 8, 9, 10, 11) on the NTU RGB+D 60 dataset: First, the proposed CAE+ achieves significant improvement (5.9%-23.9% accuracy) over three existing unsupervised methods (Id = 8, 10, 11). Compared with the shuffle&learn method (Id = 8), LongT GAN model (Id = 10), and MS$^2$L (Id = 11) that rely on challenging pretext tasks (identifying temporal order, reconstruction) and task-specific structures (deep CNNs, generative models, encoder-decoders) to learn action features, the proposed AS-CAL exploits a simpler and more flexible contrastive learning paradigm to learn effective action representations. In particular, our generic approach can be extended by different encoder structures and pretext tasks, and it can be applied to other potential skeleton-related tasks. Second, on the cross-view (C-View) testing set, the CAE+ significantly outperforms Li *et al.* (Id = 9) that applies both cross-view decoding task and reconstruction task by 10.0% accuracy improvement, which demonstrates the higher robustness of our approach against view-point changes. In addition, although our approach uses skeleton data as inputs, which are of much smaller size than depth images, it can still achieve superior performance to depth image based methods

TABLE IV

COMPARISON WITH RNN, GRU, AND LSTM MODELS ON DIFFERENT TESTING VIEWS OF UWA3D DATASET.

| | | Training Views | | V1&V2 | | V1&V3 | | V1&V4 | | V2&V3 | | V2&V4 | | V3&V4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Id | Method | Testing Views | | V3 | V4 | V2 | V4 | V2 | V3 | V1 | V4 | V1 | V3 | V1 | V2 | Average |
| | **Supervised Methods** | | | | | | | | | | | | | | | |
| 1 | RNN | | | 12.0 | 10.2 | 11.8 | 11.0 | 11.4 | 11.2 | 12.9 | 11.4 | 12.5 | 11.6 | 12.9 | 11.0 | 11.7 |
| 2 | GRU | | | 12.4 | 11.4 | 12.2 | 12.2 | 11.8 | 11.2 | 12.5 | 11.8 | 13.3 | 13.5 | 12.2 | 12.6 | 12.3 |
| 3 | LSTM | | | 12.7 | 10.2 | 11.8 | 11.0 | 12.6 | 15.5 | 12.5 | 11.4 | 12.9 | 12.4 | 12.2 | 11.4 | 12.2 |
| | **Unsupervised Methods** | | | | | | | | | | | | | | | |
| 4 | Ours (CAE) | | | 24.3 | 22.8 | 19.7 | 17.7 | 20.9 | 19.9 | 21.2 | 19.3 | 20.0 | 17.5 | 18.0 | 18.1 | 20.0 |
| 5 | Ours (CAE+) | | | **25.1** | **22.8** | **21.3** | **19.7** | **22.4** | **25.5** | **21.6** | **19.5** | **23.9** | **21.1** | **21.2** | **19.7** | **22.0** |

TABLE V

COMPARISON WITH RNN, GRU, AND LSTM MODELS ON THE SBU DATASET WITH 5-FOLD CROSS VALIDATION.

| Id | Method | Fold | | | | | |
|---|---|---|---|---|---|---|---|
| | **Supervised Methods** | 1 | 2 | 3 | 4 | 5 | Average |
| 1 | RNN | 40.0 | 42.3 | 26.8 | 27.8 | 35.4 | 34.5 |
| 2 | GRU | 40.0 | 40.4 | 28.6 | 33.3 | 40.0 | 36.5 |
| 3 | LSTM | 49.1 | 53.2 | 37.5 | 42.0 | 53.8 | 47.1 |
| | **Unsupervised Methods** | | | | | | |
| 4 | Ours (CAE) | 52.7 | 46.2 | 41.1 | 31.5 | 41.5 | 42.6 |
| 5 | Ours (CAE+) | 52.7 | 50.0 | 44.6 | 37.0 | 49.2 | 46.7 |

(Id = 1, 2, 3, 8, 9). These results indeed shows the effectiveness and efficiency of our approach.

*2) Comparison with Hand-Crafted and Supervised Methods:* The proposed action representations (CAE, CAE+) learned by AS-CAL significantly outperform existing hand-crafted methods (Id = 1-5 in Table II) on the NTU RGB+D 60 dataset. For example, our approach surpasses the representative Skeletal Quads (Id = 4) and Lie group (Id = 5) by 12.0%-23.4% on the C-View setting and 8.4%-19.9% on the C-Sub setting. Our approach is also shown to achieve comparable or even superior performance to many supervised learning methods on four datasets: **(a)** On the largest NTU RGB+D 120 dataset (see Table III), our approach performs better than the Soft RNN model and Part-Aware LSTM model by a large margin (up to 23.1% accuracy). **(b)** On the NTU RGB+D 60 dataset (see Table II), our approach obtains comparable accuracy to the deep RNN model (0.7%-2.2% accuracy improvement) and Deep RNN model (0.8% accuracy improvement on C-View). **(c)** on the SBU dataset (Table V), our approach surpasses RNN and GRU baseline models by 3.7%-17.8% accuracy on different testing folds. Despite our approach's average performance is slightly inferior to the supervised LSTM model, it attains an evidently higher performance on two of five testing folds with 3.6%-7.1% accuracy gain. Since the data size of SBU dataset (only 282 samples) is much smaller than NTU RGB+D datasets (more than 56578 samples), it is likely insufficient to train our model, *i.e.*, it performs contrastive learning with limited negative samples and small dictionary (queue), which could lead to a severe degradation of the model performance (see Sec. V-A3, Fig. 7). In addition, although our model performs action representation learning without using any skeleton label, it still achieves a comparable or even superior performance to the supervised LSTM model that leverages massive manual annotations. **(d)** On the UWA3D dataset (Table IV), our approach consistently improves the

supervised learning baselines (RNN, GRU, LSTM) by at least 6.7% accuracy on all 12 view settings, which demonstrates that our approach is more robust against view point changes than commonly used supervised methods. Since the small UWA3D dataset is challenging with limited training samples, high inter-action similarity, and frequent self-occlusions, supervised methods (Id = 1-3 in Table IV) are hard to obtain a satisfactory performance. while our approach (AS-CAL) can learn a more effective action representation under these challenges. Results (a)-(d) indicate that our unsupervised AS-CAL can achieve better performance than supervised learning baselines (RNN, GRU, LSTM) and their improved models (HBRNN, Soft RNN, Part-Aware LSTM) on both large and small datasets.

## V. DISCUSSION

To systematically evaluate the effectiveness of AS-CAL and the learned action representation, we perform ablation study on major components as well as compare different augmentation strategies, contrastive learning paradigms, and action representations.

### A. Ablation Study

In this part, we first analyze the effects of projection heads and output dimensions on our model (see Sec. V-A1). Then, we evaluate the performance of our approach under different number of layers and hidden units of the encoder (see Sec. V-A2). Finally, we thoroughly explore the influence of crucial hyper-parameters (*e.g.*, queue size $K$, momentum coefficient $m$, temperature $\tau$) in contrastive learning (see Sec. V-A3).

*1) Projection Heads and Projection Output Dimensions:* In contrastive learning, a small neural network projection head is usually used to map the learned representations to a contrastive learning space, so as to improve the quality of representations [35]. In this work, we explore two types of projection heads: (a) *Non-linear* projection head: A 2-layer multi-layer perceptron (MLP) (*i.e.*, a non-linear layer with the ReLU activation function plus a linear layer). (b) *Linear* projection head: A linear layer. We attach the non-linear or linear projection head to all encoders ($f_q$ and $f_k$), and project the action representations to a latent contrastive learning space with 64, 128, 256, 512 projection output dimensions. Note here we make a performance comparison between two action representations, namely using TAP and not using TAP (*i.e.*, use $\tilde{h}_T, \overline{h}_T$), under different projection heads (Nonlinear, Linear, None) and different projection output dimensions. As
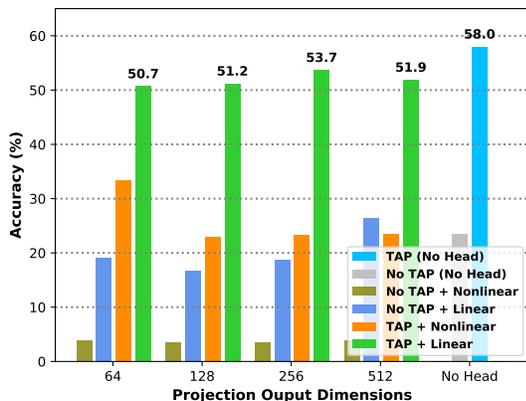
Fig. 5. Top-1 accuracy comparison on NTU RGB+D 60 (C-Sub) using different action representations (TAP / No TAP) and projection heads (Linear / Nonlinear) with various projection output dimensions for contrastive learning.

TABLE VI
TOP-1 ACCURACY (%) USING LSTM ENCODERS WITH DIFFERENT NUMBERS OF LAYERS AND HIDDEN UNITS ON NTU RGB+D 60 (C-SUB).

| Layers/Hidden Units | 64 | 128 | 192 | 256 | 320 |
|---|---|---|---|---|---|
| 1 | 12.5 | 6.5 | 49.4 | 50.1 | 47.0 |
| 2 | 8.9 | 51.3 | 55.2 | **58.0** | **58.3** |
| 3 | 1.7 | 54.6 | 55.3 | 56.3 | 56.6 |

shown in Fig. 5, we can observe some crucial results and draw conclusions as following:

**(a)** Our approach using TAP shows an evidently higher performance (20%-25% accuracy improvement) than No TAP under different settings. These results demonstrate our claim that TAP is a more effective manner to aggregate global action encoding information, which facilitate learning a better action representation. **(b)** The model attaching the linear project head significantly outperforms the one using the nonlinear head by almost double accuracy improvement, while the model without projection head (No Head) achieves the best performance over all settings. The result is different from the conclusion in [35], which claims that using non-linear head can improve the unsupervised representation learning of images. However, we argue that action representations (*i.e.*, long-term action dynamics of skeleton sequences) essentially contain more pattern information than image representations, and adding the projection head could result in action information loss due to the linear or nonlinear transformation. It can be inferred that the action information loss (note that non-linear transformation leads to more loss) makes the contrastive learning more difficult, which greatly degrades the effectiveness of final action representations. Therefore, we do NOT add the projection head to the proposed AS-CAL to keep better performance.

*2) Layers and Hidden Units of Encoder:* We take NTU RGB+D 60 (C-Sub) as an example to evaluate the effects of layers and hidden units on the performance of our approach: **(a)** As shown in Table VI and Fig. 6, using more hidden units can improve the performance under most cases (note that single-layer LSTM with a big number of hidden units degrades the performance). It can be inferred that large embedding size (*i.e.*, higher dimensional representations) is beneficial to
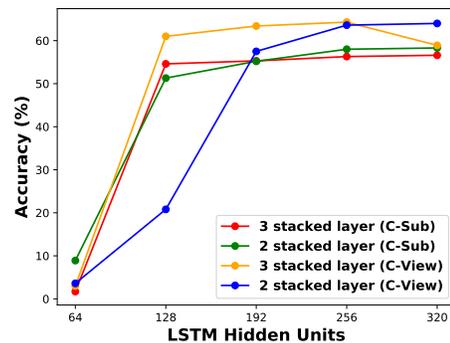


Fig. 6. Top-1 accuracy comparison using 2-layer or 3-layer LSTM encoders with different hidden units on NTU RGB+D 60 (C-Sub) and (C-View).

TABLE VII
TOP-1 ACCURACY COMPARISON OF DIFFERENT MOMENTUM COEFFICIENT $m$ ON TWO NTU RGB+D DATASETS.

| Momentum Coefficient $m$ | 0 | 0.9 | 0.99 | 0.9945 | **0.999** | 0.99945 | 0.9999 |
|---|---|---|---|---|---|---|---|
| NTU 60 (C-Sub) | 3.1 | 5.4 | 14.8 | 55.5 | **58.0** | 56.6 | 54 |
| NTU 60 (C-View) | 2.7 | 11.2 | 8.6 | 60.1 | **63.6** | 63.1 | 63.1 |
| NTU 120 (C-Sub) | 3.8 | 1.1 | 1.2 | 46.1 | **48.9** | 46.8 | 44.3 |
| NTU 120 (C-Set) | 0.8 | 0.8 | 9.6 | 48.2 | **49.7** | 49.4 | 49.1 |

aggregate more effective action features, while small number of hidden units that compress pattern information of long skeleton sequence lead to worse performance. **(b)** We observe 2-layer LSTM with 256 units achieves comparable performance to the best one (320 units). Since we expect our model to learn more compact representations with less training cost, we select 2-layer LSTM with 256 units as the query and key encoders for contrastive learning in all experiments.

*3) Queue size $K$, Momentum Coefficient $m$, Temperature $\tau$:* **(a)** As shown in Fig. 12 and Fig. 7 (note that $K$ in Fig. 12 represents the negative keys in queue (AS-CAL) and memory bank respectively), larger size of queue constantly improves the performance of our approach. It verifies the claim that more negative samples in the dictionary facilitate contrastive learning to achieve better (action) representations [34], [36]. **(b)** As presented in Sec. III-B2, the fast and stable training of contrastive learning benefits from the low and smooth update of mLSTM especially when $m = 0.999$. In Table VII, we observe that our approach also achieves the best action recognition performance on NTU RGB+D datasets with this momentum coefficient, which essentially demonstrates that a better contrastive learning is the key to achieving more effective action representations. **(c)** We evaluate the performance of our approach with different temperature $\tau$ in Table VIII, and select $\tau = 0.06$ for the proposed AS-CAL to obtain the best performance for most cases. Other datasets report similar results of (a)-(c).

*B. Comparison of Different Data Augmentations*

To systematically evaluate the effectiveness of the proposed data augmentation strategies in Sec. III-A, we take NTU RGB+D 60 dataset (C-Sub) as an example to test the performance of our approach with different compositions of augmentations. In particular, we first compare the performance of our approach between using only one augmentation strategy

TABLE VIII
TOP-1 ACCURACY COMPARISON OF DIFFERENT TEMPERATURE $\tau$ ON TWO NTU RGB+D DATASETS.

| Temperature $\tau$ | 0.03 | 0.04 | 0.05 | **0.06** | 0.07 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|---|---|
| NTU 60 (C-Sub) | 55.6 | 56.4 | 57.8 | **58.0** | 57.1 | 54.9 | 54.7 | 53.5 | 52.9 | 40.5 |
| NTU 60 (C-View) | 62.3 | 62.3 | 62.6 | **63.6** | 62.6 | 60.4 | 60.4 | 59.1 | 58.1 | 57.8 |
| NTU 120 (C-Sub) | 48.8 | 48.5 | 46.8 | **48.9** | 48.4 | 30.0 | 45.9 | 46.3 | 44.7 | 44.5 |
| NTU 120 (C-Set) | **50.8** | 49.6 | 49.8 | 49.7 | 49.0 | 47.6 | 47.9 | 46.1 | 44.5 | 44.5 |



Fig. 8. Contrastive loss curves during training solely using different augmentation strategies.



Fig. 7. Top-1 accuracy comparison of different queue sizes $K$ on two NTU RGB+D datasets.



Fig. 9. Top-1 accuracy using a single augmentation strategy on NTU RGB+D 60 (C-Sub). Note: The horizontal axis denotes different augmentation strategies.

and using the original skeleton sequence (see Fig. 9). Then, we comprehensively evaluate the effectiveness of our approach using compositions of two augmentation strategies (see Fig. 10). Finally, we empirically select several most effective augmentations to sequentially transform skeleton sequences and test the final performance on different datasets (see Fig. 11). From the results reported in Fig. 9, Fig. 10, and Fig. 11, we draw the following analysis and conclusions:

**(1)** Compared with directly using the original sequence, applying different augmentation strategies (except "CM") to AS-CAL significantly improves the accuracy by 3.4%-10.2%. As shown by Fig. 8, the contrastive loss curves of effective augmentation strategies can converge to a low loss similarly, while the "CM" curve presents a drastic fluctuation and a high contrastive loss. It suggests that a good augmentation strategy can encourage a better contrastive learning (AS-CAL), so as to achieve a more effective action representation. **(2)** Most compositions of two augmentation strategies, which transform the skeleton sequence with two different manners in order, can further boost the performance of our approach with up to 10% accuracy gain. However, double "Shear" transformations degrade the performance of propose AS-CAL, which can be inferred that drastic changes of body shape increase the difficulty to extract discriminative pattern information from skeleton sequences to recognize actions. **(3)** As reported in Fig. 10 and Fig. 11, the composition of "Reverse" and "Shear" strategies consistently achieves the best performance on action recognition when compared with other strategies. Since the sequence order typically involves the semantics of action's temporal coherence, and the shape (angle) changes of body usually contain unique pattern information of actions, this composition encourages our model to learn richer action semantics from transformed skeleton sequences for contrastive
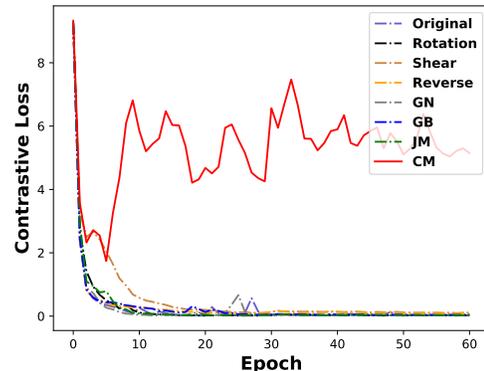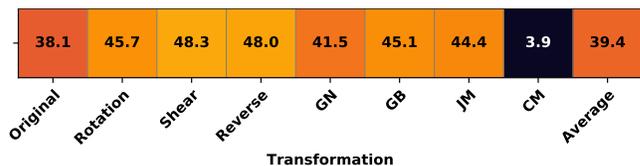
learning and action recognition. **(4)** Applying the composition of more than two augmentation strategies to AS-CAL can even perform better than using two augmentations. The composition of "Reverse" and "Shear" is the best performer in Fig. 11, and all multiple compositions based on them surpass 50% Top-1 accuracy on NTU RGB+D 60 (C-Sub), which are generally higher than compositions of two augmentations in Fig. 10.

### C. Comparison of Existing Contrastive Paradigms

In Table I, we compare the structure of proposed AS-CAL with two existing contrastive paradigms: (1) End to end paradigm using mini-batch based dictionary without momentum-based encoder [45]. (2) Memory bank paradigm with momentum update on representations of the same sample [33]. To demonstrate the effectiveness of the proposed AS-CAL, we compare the performance of three paradigms with different sizes of dictionary ($K = 128, 256, 512, 1024$). As presented in Fig. 12(a) and Fig. 12(b), AS-CAL possesses evident merits over existing contrastive paradigms in terms of action recognition performance: **(1)** The queue of AS-CAL can maintain a larger dictionary with a flexible management mechanism (illustrated in Sec. III-B4), which encourages achieving better action representation learning with an improvement of 3.0%-3.4% accuracy on the NTU RGB+D 60 dataset. Nevertheless, the large size of dictionary based on mini-batch degrades the performance of end to end model due to the increasing difficulty of large-batch optimization [46]. In addition, as analyzed in Sec. III-B2, the lack of a smooth momentum-based update of key encoder leads to a bad contrastive learning, which also prevents the end to end model (note that it directly updates encoders without using momentum) from obtaining a high accuracy. **(2)** Compared with the
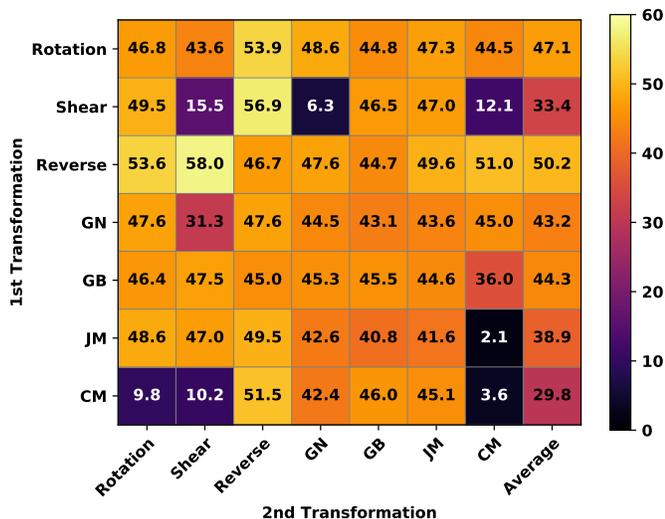
Fig. 10. Top-1 accuracy with different augmentation strategy compositions on NTU RGB+D 60 (C-Sub). Note: Every item in a row shows the accuracy of model sequentially applying two augmentations.
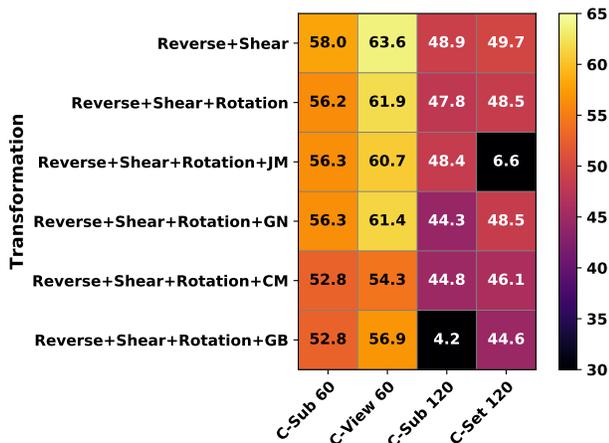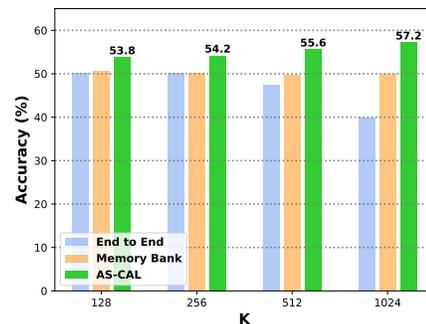


Fig. 11. Top-1 accuracy with different augmentation strategy compositions on two NTU RGB+D datasets.
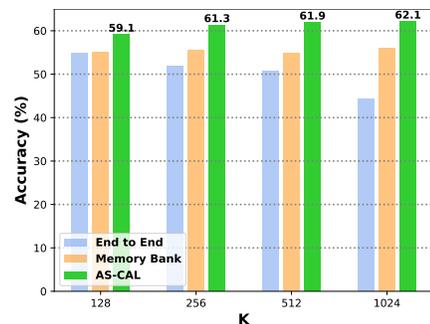
memory bank paradigm, the proposed AS-CAL benefits more from larger dictionary and can obtain higher accuracy with less memory cost (note that the memory bank needs larger memory to keep all keys in the dataset) on different settings of NTU RGB+D 60 dataset. These results also justify our claim that using queue can achieve more consistent dictionary than memory bank to achieve better action encoding.

### D. Evaluation of Different Action Representations

We comprehensively compare the performance of different action representations discussed in the paper: (1) $\tilde{h}_T$, (2) $\overline{h}_T$, (3) $k$, (4) CAE ($q$), and (5) CAE+ ($q + k$) on NTU RGB+D datasets. As reported in Table IX, the proposed CAE achieves the best performance on NTU RGB+D 120 dataset, while obtaining a comparable performance (1%-2% accuracy lower) to the CAE+ on NTU RGB+D 60 dataset. Compared with the last hidden states ($\tilde{h}_T$ and $\overline{h}_T$) that compress the temporal dynamics of a sequence [49], the action representations built by TAP, which aggregate the global action information in



(a) NTU RGB+D 60 (C-Sub)



(b) NTU RGB+D 60 (C-View)

Fig. 12. Top-1 accuracy comparison of three contrastive learning paradigms using linear evaluation for action recognition.

TABLE IX
TOP-1 ACCURACY OF DIFFERENT ACTION REPRESENTATIONS ON TWO
NTU RGB+D DATASETS.

|  | $\tilde{h}_T$ | $\overline{h}_T$ | $k$ | CAE ($q$) | CAE+ ($q + k$) |
|---|---|---|---|---|---|
| NTU 120 (C-Sub) | 2.0 | 2.1 | 48.6 | **48.9** | 48.7 |
| NTU 120 (C-Set) | 0.9 | 0.8 | 49.4 | **49.7** | 49.6 |
| NTU 60 (C-Sub) | 23.4 | 24.0 | 57.7 | 58.0 | **58.5** |
| NTU 60 (C-View) | 26.8 | 27.2 | 63.5 | 63.6 | **64.8** |

an average manner, show evidently higher effectiveness (over 30% accuracy improvement) on action recognition. Interestingly, $\tilde{h}_T$ and $\overline{h}_T$ only achieve 1%-2% Top-1 accuracy on NTU RGB+D 120 dataset. By contrast, the proposed CAE shows a stable and highly competitive performance on these larger datasets.

### E. Performance of Semi-Supervised Learning

The proposed AS-CAL could be exploited for semi-supervised learning by fine-tuning on a certain fraction (1%, 10%, 50%) of labeled data. First, we sample labeled data of NTU RGB+D datasets in a class-balanced way (i.e., around 9 (1%), 90 (10%), 450 (50%) sequences per class respectively). Then, we attach a linear classifier to the pre-trained AS-CAL model, and fine-tune the whole model with the sampled labeled data. Last, we frozen the AS-CAL model and train the linear classifier on the complete training set. Table X shows the performance of our approach under different fractions of labeled data. We discover that using the unsupervised AS-CAL (0% label) for linear evaluation can even outperform applying semi-supervised learning using labels (1% and 10% label fraction) by an evident margin (up to 13.4% Top-1

TABLE X
MODEL PERFORMANCE (TOP-1 AND TOP-5 ACCURACY) USING NO LABEL
(UNSUPERVISED LEARNING) OR FEW LABELS (SEMI-SUPERVISED
LEARNING) ON TWO NTU RGB+D DATASETS.

| Dataset/Label Fraction | Top-1 | | | | Top-5 | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | 0% | 1% | 10% | 50% | 0% | 1% | 10% | 50% |
| NTU 60 (C-Sub) | 58.0 | 47.2 | 52.2 | **61.0** | 87.4 | 81.0 | 84.0 | **88.6** |
| NTU 60 (C-View) | 63.6 | 53.5 | 57.3 | **67.3** | 91.2 | 86.5 | 88.7 | **93.0** |
| NTU 120 (C-Sub) | 48.9 | 36.0 | 42.3 | **52.6** | 78.9 | 67.1 | 74.3 | **81.3** |
| NTU 120 (C-Set) | 49.7 | 38.3 | 43.0 | **53.0** | 79.8 | 70.2 | 74.4 | **81.6** |

accuracy and Top-5 accuracy). These results suggest that the pre-trained AS-CAL is able to learn a highly effective action representation from only unlabeled data, while an insufficient fine-tuning with few labels degrades its performance. As the labeled fraction increases to 50%, our approach can benefit from semi-supervised learning using enough labels, and achieves an improvement of performance on different datasets.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we propose a generic unsupervised approach named AS-CAL to learn effective action representations from unlabeled skeleton data for action recognition. We propose to learn inherent action patterns by contrasting the similarity between augmented skeleton sequences transformed by multiple novel augmentation strategies, which enables our model to learn the invariant pattern and discriminative action features from unlabeled skeleton sequences. To facilitate better contrastive action learning, a novel momentum LSTM is proposed as the key encoder to achieve more consistent action representations. Besides, we introduce a queue to build a more consistent and memory-efficient dictionary with a flexible management of proceeding encoded keys to facilitate contrastive learning. We construct CAE as the final action representation to perform action recognition. Our approach significantly outperforms existing hand-crafted methods and unsupervised learning methods, and its performance is comparable or even superior to many supervised learning methods.

Our approach reveals considerable potentiality of unsupervised action recognition, and provides several valuable directions for future research: (1) Pretext tasks (*e.g.*, frame prediction, sequential reconstruction) could be incorporated to our approach for learning inherent high-level action semantics to improve the unsupervised contrastive learning. (2) Efficient encoders like graph convolutional network (GCN) could be exploited to learn more fine-grained spatial-temporal action features (*e.g.*, spatial co-occurrence features, inherent temporal coherence) from unlabeled skeleton data. (3) Skeleton augmentations could be further explored, while the theoretical analysis of their effects will be provided in the future work. Besides that, we expect to extend the proposed AS-CAL to multi-modal action learning for more vital vision tasks.

## REFERENCES

[1] L. Wang, D. Q. Huynh, and P. Koniusz, "A comparative review of recent kinect-based action recognition algorithms," *IEEE Trans. Image Process.*, vol. 29, pp. 15–28, 2019.

[2] B. Fernando, E. Gavves, J. M. Oramas, A. Ghodrati, and T. Tuytelaars, "Modeling video evolution for action recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5378–5387.

[3] I. Misra, C. L. Zitnick, and M. Hebert, "Shuffle and learn: unsupervised learning using temporal order verification," in *European Conference on Computer Vision*, 2016, pp. 527–544.

[4] J. Li, Y. Wong, Q. Zhao, and M. Kankanhalli, "Unsupervised learning of view-invariant action representations," in *International Conference on Neural Information Processing Systems*, 2018, pp. 1254–1264.

[5] E. Ohn-Bar and M. M. Trivedi, "Joint angles similarities and hog2 for action recognition," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, June 2013, pp. 465–470.

[6] X. Yang and Y. Tian, "Super normal vector for activity recognition using depth sequences," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2014, pp. 804–811.

[7] O. Oreifej and Z. Liu, "Hon4d: Histogram of oriented 4d normals for activity recognition from depth sequences," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2013, pp. 716–723.

[8] S. Xu, H. Rao, H. Peng, X. Jiang, Y. Guo, X. Hu, and B. Hu, "Attention based multi-level co-occurrence graph convolutional lstm for 3d action recognition," *IEEE Internet Things J.*, pp. 1–1, 2020.

[9] C. Si, W. Chen, W. Wang, L. Wang, and T. Tan, "An attention enhanced graph convolutional lstm network for skeleton-based action recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1227–1236.

[10] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018, pp. 7444–7452.

[11] N. Zheng, J. Wen, R. Liu, L. Long, J. Dai, and Z. Gong, "Unsupervised representation learning with long-term dynamics for skeleton based action recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018, pp. 2644–2651.

[12] Z. Ghahramani, "Unsupervised learning," in *Summer School on Machine Learning*, 2003, pp. 72–112.

[13] L. Devillers, L. Vidrascu, and L. Lamel, "Challenges in real-life emotion annotation and machine learning based detection," *Neural Netw.*, vol. 18, no. 4, pp. 407–422, 2005.

[14] T. Nguyen-Phuoc, C. Li, L. Theis, C. Richardt, and Y.-L. Yang, "Hologan: Unsupervised learning of 3d representations from natural images," in *IEEE International Conference on Computer Vision*, 2019, pp. 7588–7597.

[15] K. Su, X. Liu, and E. Shlizerman, "Predict & cluster: Unsupervised skeleton based action recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9631–9640.

[16] L. Lin, S. Song, W. Yang, and J. Liu, "Ms2l: Multi-task self-supervised learning for skeleton based action recognition," in *ACM International Conference on Multimedia*, 2020, pp. 2490–2498.

[17] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *International Conference on Neural Information Processing Systems*, 2012, pp. 1097–1105.

[18] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[19] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *International Conference on Neural Information Processing Systems 27*, 2014, pp. 2672–2680.

[20] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," *arXiv preprint arXiv:1906.05849*, 2019.

[21] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," in *International Conference on Learning Representations*, 2019.

[22] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.

[23] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[24] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+ d: A large scale dataset for 3d human activity analysis," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1010–1019.

[25] J. Liu, A. Shahroudy, M. Perez, G. Wang, L. Y. Duan, and A. C. Kot, "Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 10, pp. 2684–2701, 2020.

[26] K. Yun, J. Honorio, D. Chattopadhyay, T. L. Berg, and D. Samaras, "Two-person interaction detection using body-pose features and multiple instance learning," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2012, pp. 28–35.

[27] H. Rahmani, A. Mahmood, D. Q. Huynh, and A. Mian, "Hopc: Histogram of oriented principal components of 3d pointclouds for action recognition," in *European Conference on Computer Vision*, 2014, pp. 742–757.

[28] G. Evangelidis, G. Singh, and R. Horaud, "Skeletal quads: Human action recognition using joint quadruples," in *International Conference on Pattern Recognition*, 2014, pp. 4513–4518.

[29] R. Vemulapalli, F. Arrate, and R. Chellappa, "Human action recognition by representing 3d skeletons as points in a lie group," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 588–595.

[30] N. Srivastava, E. Mansimov, and R. Salakhudinov, "Unsupervised learning of video representations using lstms," in *International Conference on Machine Learning*, 2015, pp. 843–852.

[31] B. Su, J. Zhou, X. Ding, and Y. Wu, "Unsupervised hierarchical dynamic parsing and encoding for action recognition," *IEEE Trans. Image Process.*, vol. 26, no. 12, pp. 5784–5799, 2017.

[32] U. Ahsan, R. Madhok, and I. Essa, "Video jigsaw: Unsupervised learning of spatiotemporal context for video action recognition," in *IEEE Winter Conference on Applications of Computer Vision*, 2019, pp. 179–189.

[33] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3733–3742.

[34] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2020, pp. 9729–9738.

[35] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International Conference on Machine Learning*, 2020.

[36] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.

[37] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *International Conference on Machine Learning*, 2008, pp. 1096–1103.

[38] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *European Conference on Computer Vision*, 2016, pp. 649–666.

[39] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 2, 2006, pp. 1735–1742.

[40] 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.

[41] "Unsupervised learning of visual features by contrasting cluster assignments," in *Caron, Mathilde and Misra, Ishan and Mairal, Julien and Goyal, Priya and Bojanowski, Piotr and Joulin, Armand*, vol. 33, 2020.

[42] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton, "Big self-supervised models are strong semi-supervised learners," *arXiv preprint arXiv:2006.10029*, 2020.

[43] X. Chen and K. He, "Exploring simple siamese representation learning," *arXiv preprint arXiv:2011.10566*, 2020.

[44] H. Wang and L. Wang, "Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 499–508.

[45] M. Ye, X. Zhang, P. C. Yuen, and S.-F. Chang, "Unsupervised embedding learning via invariant and spreading instance feature," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6210–6219.

[46] P. Goyal, P. Dollár, R. B. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch sgd: Training imagenet in 1 hour," *arXiv preprint arXiv:1706.02677*, 2017.

[47] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1110–1118.

[48] J.-F. Hu, W.-S. Zheng, L. Ma, G. Wang, J. Lai, and J. Zhang, "Early action prediction by soft regression," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 11, pp. 2568–2583, 2018.

[49] J. Weston, S. Chopra, and A. Bordes, "Memory networks," in *International Conference on Learning Representations*, 2015.