

---

# Unsupervised Motion Representation Learning with Capsule Autoencoders

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose the Motion Capsule Autoencoder (MCAE), which addresses a key  
2 challenge in the unsupervised learning of motion representations: transforma-  
3 tion invariance. MCAE models motion in a two-level hierarchy. In the lower  
4 level, a spatio-temporal motion signal is divided into short, local, and semantic-  
5 agnostic *snippets*. In the higher level, the snippets are aggregated to form full-  
6 length semantic-aware *segments*. For both levels, we represent motion with a  
7 set of learned transformation invariant templates and the corresponding geometric  
8 transformations by using capsule autoencoders of a novel design. This leads to  
9 a robust and efficient encoding of viewpoint changes. MCAE is evaluated on a  
10 novel Trajectory20 motion dataset and various real-world skeleton-based human  
11 action datasets. Notably, it achieves better results than baselines on Trajectory20  
12 with considerably fewer parameters and state-of-the-art performance on the unsu-  
13 pervised skeleton-based action recognition task.

## 14 1 Introduction

15 Real-world movements contain a plethora of information beyond the literal sense of moving. For  
16 example, honeybees “dance” to communicate the location of a foraging site and human gait alone  
17 can reveal activities and identities [6]. Understanding these movements is vital for an artificial in-  
18 telligent agent to comprehend and interact with the ever-changing world. Studies on social behavior  
19 analysis [4, 5], action recognition [50, 55], and video summarizing [51] have also acknowledged the  
20 importance of movement.

21 A key step towards understanding movements is to analyze its patterns. However, learning motion  
22 pattern representations is non-trivial due to (1) the curse of dimensionality from input data, (2) diffi-  
23 culties in modeling long-term dependencies in motion sequences, (3) high intra-class variation as a  
24 result of subject or viewpoint change, and (4) insufficient data annotation. The first two challenges  
25 have been ameliorated by the advances in keypoint detection methods and spatial-temporal feature  
26 extractors [32, 37, 43]. The third and the fourth nonetheless remain hurdles and call for unsupervised  
27 transformation-invariant motion models.

28 Inspired by the viewpoint-invariant capsule-based representation for images [7, 11], we exploit cap-  
29 sule network and introduce the Motion Capsule Autoencoder (MCAE), an unsupervised capsule  
30 framework that learns the transformation-invariant motion representation for keypoints. MCAE  
31 models motion signals in a two-level snippet-segment hierarchy. At the lower level, motion signal  
32 is encoded as the snippet capsules (SniCap) that describe movement in narrow time spans (i.e. snip-  
33 pets). The motion snippets are then temporally concatenated to form movement with wider time  
34 spans (i.e. segments). At the higher level, the segment capsules (SegCap) maintain a set of segment  
35 templates (i.e. learned canonical motion patterns) and transform them to reconstruct the input seg-  
36 ment motion. Both SniCaps and SegCaps learn transformation-invariant motion representation in

37 their temporal receptive field. The SegCaps, which are built upon SniCaps, produce a high-level ab-  
38 straction of the given motion signal, where the semantics are represented by transformation-invariant  
39 capsule activations.

40 The contributions of this work are as follows:

- 41 • We propose MCAE, an unsupervised capsule framework that learns a transformation-  
42 invariant, discriminative, and compact representation of motion signals. Two motion cap-  
43 sules are designed to generate representation at different abstraction levels. The lower-level  
44 representation captures the local short-time movements, which are aggregated into higher-  
45 level representation that is discriminative for motion with wider time spans.
- 46 • We propose Trajectory20, a novel and challenging synthetic dataset with a wide class of  
47 motion patterns and controllable intra-class variations.
- 48 • Extensive experiments on both the synthetic and real-world skeleton human action datasets  
49 show the efficacy of MCAE. In addition, we perform an ablation study to examine the effect  
50 of different regularizers and some key hyperparameters of the proposed MCAE.

## 51 2 Related Works

52 **Motion Representation** A variety of methods have been proposed to learn (mostly human) mo-  
53 tion representation from video frames [1, 17, 22, 46], depth maps [8, 16, 21, 31, 41], or key-  
54 points/skeletons [2, 15, 18, 20, 23, 27, 35, 45, 48, 49]. Earlier works use handcrafted features  
55 like Fourier coefficients [41], dense trajectory features [24, 40], and Lie group representations [38].  
56 Some works use canonical human pose [26] or view-invariant short tracklets to learn robust feature  
57 for recognition [13]. The development of deep learning brings the usage of convolution networks  
58 (ConvNet) and recurrent networks for motion representation. Simonyan *et al.* [33] proposes a two-  
59 stream ConvNet which combines video frame with optical flow. C3D [37] proposes to use 3D  
60 convolution on the spatial-temporal cubes. Srivastava *et al.* [34] uses an LSTM-based encoder to  
61 map input frames to a fixed-length vector and apply task-dependent decoders for applications such  
62 as frame reconstruction and frame prediction. The combined use of convolution module and LSTM  
63 has also been proved effective in [1, 32, 44].

64 A series of works [9, 12, 14, 39, 42] have been proposed to address the problem of learning  
65 viewpoint-invariant motion representation from videos or keypoint sequences. Most of these meth-  
66 ods rely on multi-modality input of RGB frames, depth maps or keypoint trajectories. Some other  
67 works [18, 27, 35, 54] focus on the unsupervised learning of keypoint/skeleton motion. In these  
68 works, LSTM is widely used for motion modelling which generally results in a heavy memory  
69 footprint.

70 Different from prior works, MCAE takes only keypoint motion as input and uses LSTM to model  
71 snippet-segment relations only. This results in a lightweight yet capable model.

72 **Capsule Network** MCAE is closely related to the Capsule Network [28], which is designed to  
73 represent objects in images using automatically discovered constituent parts and poses. The explicit  
74 modeling of poses helps learning viewpoint-invariant visual features that are more compact and  
75 flexible than traditional ConvNets. Kosiorek *et al.* [11] proposed the unsupervised stacked capsule  
76 autoencoder (SCAE), which learns view-invariant representation for images. More recently, capsule  
77 network has also been applied to point cloud processing [52, 53] for 3D object classification and  
78 reconstruction.

79 Despite the success of capsule networks in various vision tasks, the study of capsule networks on  
80 motion representation is scarce. VideoCapsuleNet [3] proposes to generalize capsule networks from  
81 2D to 3D for action detection in videos. Yu *et al.* [47] proposed a limited study on supervised  
82 skeleton-based action recognition using Capsule Network. Sankisa *et al.* [30] proposed to use Cap-  
83 sule Network for error concealment in videos.

84 Different from these works, MCAE performs unsupervised learning of motion represented as co-  
85 ordinates rather than pixels. It aims at learning an appearance-agnostic transformation-invariant  
86 motion representation.

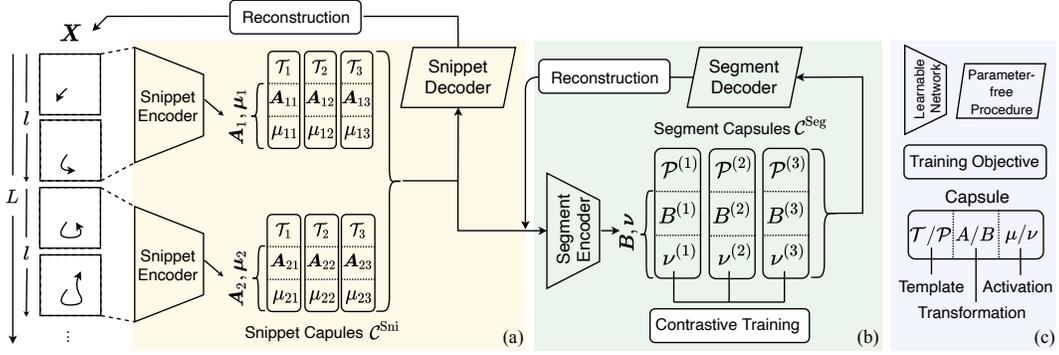


Figure 1: Overview of MCAE (best viewed in color). (a) The Snippet Autoencoder, which learns the semantic-agnostic short-time representation (snippet capsules) by reconstructing the input signal  $\mathbf{X}$ . (b) The Segment Autoencoder, which learns the semantic-aware long-time representation (segment capsules) by aggregating and reconstructing snippet capsule parameters. The activation values in segment capsules are used as semantic information for self-supervised contrastive training. (c) Meanings for different shapes and variables.

### 87 3 Methodology

88 We consider a single point<sup>1</sup> in  $d$ -dimension space. The motion of the point, i.e. a trajectory, is  
 89 described by  $\mathbf{X} = \{\mathbf{x}_i | i = 1, \dots, L\}$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  is the coordinates at time  $i$ . Semantically,  $\mathbf{X}$   
 90 belongs to a motion pattern, subject to an arbitrary and unknown geometric transformation. Given  
 91 sufficient samples of  $\mathbf{X}$ , we aim to learn a discriminative (in particular, transformation-invariant)  
 92 representation for those motion samples without supervision.

#### 93 3.1 Framework Overview

94 We solve this problem in two steps, namely *snippet learning* and *segment learning*. Snippets and  
 95 segments correspond to the lower and higher levels in the hierarchy of how MCAE views the motion  
 96 signal. Both snippets and segments are temporally consecutive subsets of  $\mathbf{X}$ , but snippets have a  
 97 shorter time span than segments. In the snippet learning step, the input  $\mathbf{X}$  is first divided into  $L/l$   
 98 temporally non-overlapping snippets, where  $l$  is the length of snippets. Each of these snippets will  
 99 be mapped into a semantic-agnostic representation by the **Snippet Autoencoder**. In the segment  
 100 learning step, the snippet representations are combined and fed into the **Segment Autoencoder**,  
 101 where the full motion is represented as a weighted mixture of the transformed canonical represen-  
 102 tations. The segment activations are used as the motion representation for downstream tasks. An  
 103 overview of the framework is shown in Fig. 1. In the following section, we delineate the details for  
 104 each module and explain the training procedure.

#### 105 3.2 Snippet Autoencoder

106 To encode the snippets' motion variation, we propose the Snippet Capsule (SniCap), which we  
 107 denote as  $\mathcal{C}^{\text{Sni}}$ . SniCap is parameterized as  $\mathcal{C}^{\text{Sni}} = \{\mathcal{T}, \mathbf{A}, \mu\}$ , where  $\mathcal{T}$ ,  $\mathbf{A}$ , and  $\mu$  are the **snippet**  
 108 **template**, **snippet transformation parameter**, and **snippet activation**, respectively. The snippet  
 109 template  $\mathcal{T} = \{t_i | t_i \in \mathbb{R}^d, i = 1, \dots, l\}$  describes a motion template of length  $l$  and is the identity  
 110 information of a SniCap.  $\mathbf{A}$  and  $\mu$  depend on the input snippet. The transformation parameter  
 111  $\mathbf{A} \in \mathbb{R}^{(d+1) \times (d+1)}$  describes the geometric relation between the input snippet and the snippet template.  
 112 The snippet activation  $\mu \in [0, 1]$  denotes whether the snippet template is used, i.e. activated, to  
 113 represent the input snippet.

114 **Snippet Encoding/Decoding** For a given snippet  $\mathbf{x}_{i:i+l}$ , the snippet module performs the follow-  
 115 ing steps: (1) encode motion properties with Snippet Encoder into SniCaps, and (2) decode SniCaps  
 116 to reconstruct the original  $\mathbf{x}_{i:i+l}$ . For the encoding step, a 1D-ConvNet  $f_{\text{CONV}}$  is used to extract

<sup>1</sup>We show a way to generalize MCAE to multi-point systems in Section 4.2

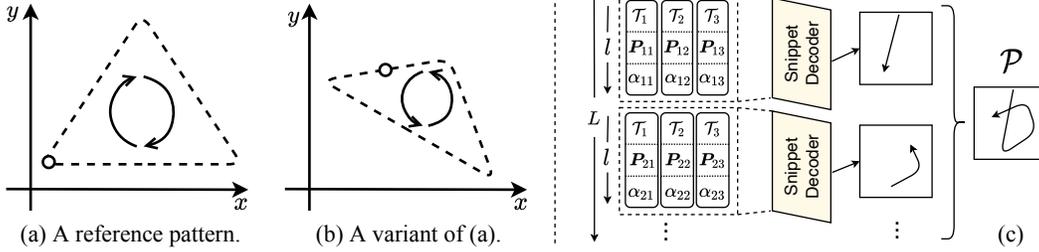


Figure 2: (a) and (b) show a reference motion pattern and a variant of it. The circle and the arrow shows the start and the direction of motion respectively. (c) Interpretation of a segment template  $\mathcal{P}$ .  $\mathcal{P}$  is functionally the same as  $S$  snippet parameters  $(\mathbf{A}, \mu)$ . When combined with  $\mathcal{T}$ , it can be decoded into an  $L$ -long sequence. The segment autoencoder maintains multiple segment templates, which can be transformed and mixed to reconstruct the input snippet parameters.

117 the motion information from  $\mathbf{x}_{i:i+l}$  and predict SniCap parameters, i.e.  $\{(\mathbf{A}_j, \mu_j) | j = 1, \dots, N\} =$   
 118  $f_{\text{CONV}}(\mathbf{x}_{i:i+l})$  where  $N$  is the number of SniCaps. For the decoding step, we first apply the transfor-  
 119 mation  $\mathbf{A}$  to the snippet templates as

$$\begin{pmatrix} \hat{\mathbf{t}}_{ij} \\ 1 \end{pmatrix} = \mathbf{A}_i \begin{pmatrix} \mathbf{t}_j \\ 1 \end{pmatrix}, \quad i = 1, \dots, N, \quad j = 1, \dots, l. \quad (1)$$

120 Then, the transformed templates from different SniCaps are mixed, according to their activations,  
 121 and the corresponding reconstructed input is

$$\hat{\mathbf{x}}_j = \sum_{i=1}^N \mu_i \hat{\mathbf{t}}_{ij}, \quad j = 1, \dots, l, \quad (2)$$

122 where  $\hat{\mathbf{t}}_{ij}$  indicates the transformed coordinate of the  $i^{\text{th}}$  SniCap at  $j^{\text{th}}$  time step.

### 123 3.3 Segment Autoencoder

124 The motion information encoded in SniCaps is agnostic to the segment level motion patterns. This  
 125 makes it less biased towards the training data domain. However, its utility on high-level applications,  
 126 such as activity analysis or motion classification, is greatly undermined. For example, consider  
 127 Fig. 2(a) as a reference “triangle” trajectory. Fig. 2(b) illustrates a possible variation. Since the two  
 128 trajectories differ a lot in their local movement, they could be considered as different classes without  
 129 transformation-invariant information from the full trajectory.

130 Hence, we introduce a segment encoder to gain a holistic understanding of motion and encapsu-  
 131 late such information in the segment capsules (SegCap). A SegCap is parameterized as  $\mathcal{C}^{\text{Seg}} =$   
 132  $\{\mathcal{P}, \mathbf{B}, \nu\}$ , where  $\mathcal{P}$ ,  $\mathbf{B}$ , and  $\nu$  are the **segment template**, **segment transformation parameter**, and  
 133 **segment activation**, respectively. The segment template  $\mathcal{P}$  is fixed for a SegCap w.r.t the training  
 134 domain. It describes a set of canonical motion patterns in terms of all the snippet templates and is  
 135 defined as  $\mathcal{P} = \{(\mathbf{P}_i, \alpha_i) | i = 1, \dots, S\}$ , where  $\mathbf{P}_i \in \mathbb{R}^{N \times (d+1) \times (d+1)}$  and  $\alpha_i \in \mathbb{R}^N$ .  $S = L/l$  is the  
 136 number of snippets. Each  $\mathbf{P}_{ij} \in \mathbb{R}^{(d+1) \times (d+1)}$  describes how the  $j^{\text{th}}$  snippet template in SniCaps  
 137 is aligned to form the motion pattern. The weight  $\alpha_{ij}$  controls the importance of  $j^{\text{th}}$  snippet tem-  
 138 plate in the  $i^{\text{th}}$  snippet. In other words, a  $(\mathbf{P}_i, \alpha_i)$  describes how the  $N$  snippet templates are used  
 139 to construct an  $l$ -long snippet and a SegCap requires  $S$  such parameters to describe a full  $L$ -long  
 140 sequence. Fig. 2(c) illustrates the interpretation of  $\mathcal{P}$ . Both  $\mathbf{B}$  and  $\nu$  are dependent on the input.  
 141  $\mathbf{B} \in \mathbb{R}^{(d+1) \times (d+1)}$  is a transformation on  $\mathcal{P}$ , and  $\nu \in [0, 1]$  is the activation of the segment template.

142 **Segment Encoding/Decoding** Assume we have  $M$  SegCaps with which we hope to reconstruct  
 143 the low-level motion encoded in the SniCap parameters. This is equivalent to reconstructing all  
 144 the data-dependent SniCap parameters  $[\mathcal{C}_1^{\text{Sni}}, \dots, \mathcal{C}_S^{\text{Sni}}]$ , where  $\mathcal{C}_i^{\text{Sni}} = \{(\mathbf{A}_j, \mu_j) | j = 1, \dots, N\}$  is  
 145 the set of SniCap parameters for the  $i^{\text{th}}$  snippet. To obtain the SegCap parameters, we encode the  
 146  $S$ -long sequence of SniCap parameters with an LSTM model  $f_{\text{LSTM}}$  shared by all SegCaps, and  $M$

147 fully-connected layers (one for each SegCap) to produce  $\{\mathbf{B}, \nu\}$ . Formally,

$$\begin{aligned} \mathbf{h} &= f_{\text{LSTM}}\left([\mathcal{C}_1^{\text{Sni}}, \dots, \mathcal{C}_S^{\text{Sni}}]\right), \\ \{\mathbf{B}^{(k)}, \nu^{(k)}\} &= f_{\text{FC}}^{(k)}(\mathbf{T}, \mathbf{h}), \quad k = 1, \dots, M, \end{aligned} \quad (3)$$

148 where superscript  $(k)$  refers to the  $k^{\text{th}}$  SegCap. The transformation and activation parameters are  
149 then applied to  $\mathcal{P}$  to reconstruct snippet parameters

$$\begin{aligned} \hat{\mathbf{P}}_{ij}^{(k)} &= \mathbf{B}^{(k)} \times \mathbf{P}_{ij}^{(k)}, \quad i = 1, \dots, S, \quad j = 1, \dots, N, \quad k = 1, \dots, M, \\ \hat{\mathcal{C}}_i^{\text{Sni}} &= (\hat{\mathbf{A}}_i, \hat{\boldsymbol{\mu}}_i) = \left( \sum_k \nu^{(k)} \hat{\mathbf{P}}_i^{(k)}, \sum_k \nu^{(k)} \boldsymbol{\alpha}_i^{(k)} \right), \quad i = 1, \dots, S, \end{aligned} \quad (4)$$

150 where  $\hat{\mathbf{A}}_i \in \mathbb{R}^{N \times (d+1) \times (d+1)}$  and  $\hat{\boldsymbol{\mu}}_i \in \mathbb{R}^N$  are the reconstructed snippet transformation and acti-  
151 vation of the snippet templates for the  $i^{\text{th}}$  snippet. Note that  $S = L/l$ , which means  $f_{\text{LSTM}}$  can have  
152 a much smaller footprint than a recurrent network that handles the whole  $L$ -long sequence.

153 The above formulation enables SegCap to learn a transformation-invariant representation of motion.  
154 Intuitively,  $\mathcal{P}$  describes snippet-segment relation, and  $\mathbf{B}$  can be regarded as the spatial relation  
155 between a segment template pattern and the observed trajectory. The segment activation  $\nu \in \mathbb{R}^M$   
156 reveals the semantics of the input trajectory and can be used for self-supervised training.

### 157 3.4 Training

158 As delineated in Section 3.2 and 3.3, SniCap and SegCap play different roles by capturing infor-  
159 mation at two different abstraction levels. SniCap focuses on short-time motion while SegCap is  
160 defined upon SniCap to model long-time semantic information. Hence, the two autoencoders are  
161 trained using different objective functions.

162 The only objective of the snippet autoencoder is to faithfully reconstruct the original input. There-  
163 fore, for a training sample  $\mathbf{X} = \{\mathbf{x}_i | i = 1, \dots, L\}$ , we use a self-supervised reconstruction loss:

$$\mathcal{L}_{\text{Rec}}^{\text{Sni}} = \sum_{i=1}^L \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2^2, \quad (5)$$

164 where  $\hat{\mathbf{x}}_i$  denotes the reconstructed coordinate following Equation (2).

165 The segment autoencoder’s primary goal is to reconstruct the input SniCap parameters, hence the  
166 reconstruction loss

$$\mathcal{L}_{\text{Rec}}^{\text{Seg}} = \sum_{i=1}^S \|\hat{\mathbf{A}}_i - \mathbf{A}_i\|_2^2 + \|\hat{\boldsymbol{\mu}}_i - \boldsymbol{\mu}_i\|_2^2. \quad (6)$$

167 Furthermore, we use unsupervised contrastive training to learn semantic meaningful activations  $\nu$ .  
168 For a batch of  $B$  samples, the contrastive loss is

$$\mathcal{L}_{\text{Con}}^{\text{Seg}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\text{cossim}(\boldsymbol{\nu}'_i, \boldsymbol{\nu}''_i)/\tau)}{\sum_j \exp(\text{cossim}(\boldsymbol{\nu}'_i, \boldsymbol{\nu}''_j)/\tau)}, \quad (7)$$

169 where  $\tau = 0.1$  is the temperature used for all experiments,  $\boldsymbol{\nu}'_i$  and  $\boldsymbol{\nu}''_i$  is the segment activation of  
170 sample  $\mathbf{X}'_i$  and  $\mathbf{X}''_i$ , respectively. Here,  $\mathbf{X}'_i$  and  $\mathbf{X}''_i$  are the spatial-temporally disturbed version of  
171  $\mathbf{X}_i$ . The disturbance is dataset-dependent and will be discussed in supplementary materials.

172 In addition to the above loss terms, we impose two regularizers: a smoothness constraint on re-  
173 constructed sequence, and a sparsity regularization on the segment activations

$$\mathcal{L}_{\text{Smt}}^{\text{Reg}} = \sum_{i=2}^L \|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{i-1}\|_2^2, \quad \mathcal{L}_{\text{Sps}}^{\text{Reg}} = \frac{1}{B} \sum_{i=1}^B \|\boldsymbol{\nu}_i\|_2^2. \quad (8)$$

174 The final training objective is:

$$\mathcal{L} = \lambda^{\text{Sni}} \mathcal{L}_{\text{Rec}}^{\text{Sni}} + \lambda^{\text{Seg}} \mathcal{L}_{\text{Rec}}^{\text{Seg}} + \mathcal{L}_{\text{Con}}^{\text{Seg}} + 0.5 \mathcal{L}_{\text{Smt}}^{\text{Reg}} + 0.05 \mathcal{L}_{\text{Sps}}^{\text{Reg}}, \quad (9)$$

175 where  $\lambda^{\text{Sni}}$  and  $\lambda^{\text{Seg}}$  are empirically determined.

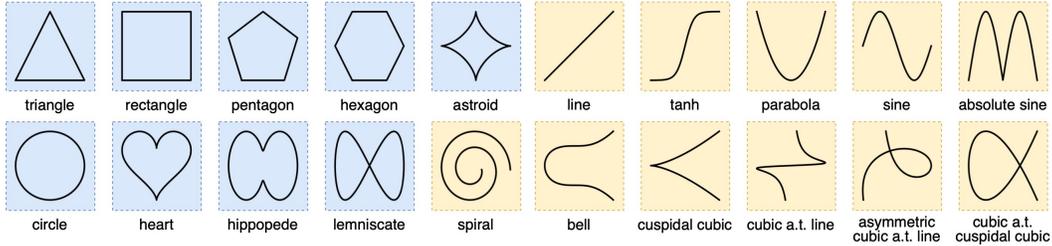


Figure 3: The 20 motion patterns in the Trajectory20 (T20) dataset. a.t. is short for “asymptotic to”.

## 176 4 Experiments

177 In this section, we first assess the proposed MCAE on a synthetic motion dataset to show its ability in  
 178 learning transformation-invariant robust representations. Then, we generalize MCAE to multi-point  
 179 systems and show its efficacy in real-world skeleton-based human action datasets. We report the  
 180 mean accuracy and standard error based on three runs with random initialization. The experiments  
 181 are run on an NVIDIA Titan V GPU, where we use a batch size of 64, and the Adam [10] optimizer  
 182 with a learning rate of  $10^{-3}$ . Please refer to the supplementary material for details.

### 183 4.1 Learning from Synthesized Motion

184 **The Trajectory20 Dataset** While datasets like movingMNIST [34] have been commonly used in  
 185 the motion representation learning literature, it is innately *linear* and has limited motion variations.  
 186 Moreover, its prediction-oriented setting makes it difficult to examine the motion category of each  
 187 trajectory. In this paper, we introduce the Trajectory20 (T20), a synthetic trajectory dataset based on  
 188 20 distinct motion templates (as shown in Fig. 3). Each sample in T20 is a 32-step-long sequence  
 189 of coordinates in  $[-1, 1]^2$ . In the data generating process, a motion template is picked randomly and  
 190 is randomly rotated, scaled, and translated to a random position to produce a trajectory. A closed  
 191 trajectory (marked blue in Fig. 3) starts at a random point on the trajectory and end at the same  
 192 point, whereas an open trajectory (marked yellow in Fig. 3) starts at a random end’s vicinity. The  
 193 randomized generating process ensures the trajectories are controllably diverse in scale, rotation,  
 194 and position. The training data is generated on-the-fly and a fixed test set of 10,000 samples is used  
 195 for evaluation. Examples of T20 are shown in the supplementary material.

196 **Ablation Study** We perform an ablation study of  
 197 MCAE on T20 to examine the effect of different  
 198 regularizers and three key hyperparameters: snippet  
 199 length  $l$ , the numbers of SniCap (#Sni) and SegCap  
 200 (#Seg). The result is shown in Table 1. The length  
 201 of snippet  $l$  plays a vital role in learning a useful  
 202 representation. A very small  $l$  results in a narrow  
 203 receptive field for snippet capsules, which makes it  
 204 less useful for inferring semantics of the whole se-  
 205 quence. At the other end, a large  $l$  makes snippets  
 206 challenging to reconstruct. The numbers of Sni-  
 207 Cap and SegCap also have major effect on the out-  
 208 come. Too few SniCaps makes it difficult to recon-  
 209 struct the input motion signal. Too few SegCaps un-  
 210 dermines the expressiveness of the segment autoen-  
 211 coder. Too many SniCaps could cause difficulty in  
 212 learning proper alignments between SegCaps and SniCaps. Both degrade the quality of the learned  
 213 features. Moreover, increasing #Seg from 80 to 128 does not bring further improvements. As the re-  
 214 sult shows,  $(l, \text{\#Sni}, \text{\#Seg}) = (8, 8, 80)$  performs well and we will use it in all experiments below. As  
 215 for the regularizers, while both regularizers improve the performance, the sparsity regulation ( $\mathcal{L}_{\text{Sps}}^{\text{Reg}}$ )  
 216 on segment activation is more helpful for learning discriminative features.

Table 1: Ablation study on T20.

Reg.	$l$	#Sni	#Seg	Acc. (%)
	8	8	80	$69.30 \pm 0.76$
	4	8	80	$41.01 \pm 8.81$
	16	8	80	$45.83 \pm 8.36$
Full	8	2	80	$64.02 \pm 2.10$
	8	4	80	$68.17 \pm 0.36$
	8	16	80	$48.11 \pm 1.60$
	8	8	32	$42.36 \pm 3.15$
	8	8	64	$63.94 \pm 1.41$
	8	8	128	$69.44 \pm 1.69$
w/o $\mathcal{L}_{\text{Smt}}^{\text{Reg}}$	8	8	80	$67.60 \pm 1.69$
w/o $\mathcal{L}_{\text{Sps}}^{\text{Reg}}$	8	8	80	$65.92 \pm 1.63$

217 **Motion Classification** We compare MCAE with the following baseline models, namely KMeans,  
 218 DTW-KMeans,  $k$ -Shape [25], LSTM and 1D-Conv<sup>2</sup>. KMeans, DTW-KMeans, and  $k$ -Shape are  
 219 parameter-free time series clustering algorithms. Briefly, KMeans uses Euclidean distance to mea-  
 220 sure the similarity between signals. DTW-KMeans normalizes input signals using dynamic time  
 221 warping [29], and performs KMeans on the normalized signals.  $k$ -Shape uses cross-correlation  
 222 based distance measure to cluster time series. We use the implementation by `tslearn` [36] for the  
 223 three clustering methods. LSTM, 1D-Conv, and MCAE are used as backbone networks, which take  
 224 the raw coordinate sequence as input and output a feature vector of a pre-defined dimension. The  
 225 feature vector is used for contrastive learning following Equation (7). Upon each model we attach  
 226 an auxiliary linear classifier (i.e. a single layer perceptron), which is trained with labels but the gra-  
 227 dient is blocked from back-propagating into the backbone. The corresponding accuracy reflects the  
 228 quality of the learned representation.

229 For LSTM and 1D-Conv backbone, dif-  
 230 ferent numbers of hidden units/channels  
 231 have been explored (shown as Hidden  
 232 Param. in Table 2), which has resulted  
 233 in different model sizes (measured by  
 234 #Param. in Table 2).

235 As shown in Table 2, since the spatial  
 236 variance (e.g. viewpoint changes) within  
 237 motion signal cannot be directly cap-  
 238 tured by temporal warping/correlation,  
 239 all the three parameter-free cluster-  
 240 ing methods perform poorly on T20.  
 241 On the other hand, with considerably  
 242 fewer parameters, MCAE outperforms  
 243 LSTM and 1D-CNN by a large margin.  
 244 This provides quantitative evidence that  
 245 MCAE can capture the transformation-  
 246 invariant semantic information more ef-  
 247 ficiently than the compared baselines.

Table 2: Unsupervised learning performance of MCAE and baselines on T20.

	Hidden Param.	#Param.	Acc. (%)
KMeans	–	–	8.57 ± 0.04
DTW-KMeans	–	–	9.12 ± 0.20
$k$ -Shape [25]	–	–	12.94 ± 0.34
	128	600k	29.17 ± 2.45
	256	669k	40.03 ± 0.57
LSTM	512	805k	45.59 ± 1.37
	1,024	1,078k	53.47 ± 1.52
	2,048	1,625k	54.32 ± 0.55
	128	588k	44.78 ± 0.57
	256	787k	53.69 ± 0.53
1D-Conv	512	1,185k	57.57 ± 0.56
	1,024	1,982k	57.58 ± 0.08
	(#Sni, #Seg)	#Param.	Acc. (%)
MCAE	(8, 80)	<b>277k</b>	<b>69.30 ± 0.76</b>

## 248 4.2 Generalizing to Multiple Points

249 The MCAE running on T20 dataset handles a single moving point while most real-world problems  
 250 involve multiple points. This section presents a naive (yet effective) extension of MCAE, which  
 251 we name MCAE-MP, to enable processing the motion of multi-point systems. Such motion can  
 252 be described as  $\mathcal{X} = \{\mathcal{X}_i | i = 1, \dots, K\}$ , where  $K$  is the number of moving points. The extension  
 253 works as follows:

- 254 1. The  $K$  moving points are processed separately by an MCAE. This results in  $K$  segment  
 255 activation vectors  $\{\nu_i, |i = 1, \dots, K\}$ .
- 256 2. The  $K$  activation vectors are concatenated into a single representation  $\nu \in \mathbb{R}^{KM}$ , which is  
 257 used for unsupervised learning following Equation (9).

258 **Skeleton-based Human Action Recognition** We apply MCAE-MP to solve the unsupervised  
 259 skeleton-based action recognition problem, where a human skeleton is a system consisting of mul-  
 260 tiple moving joints (points). Three widely-used datasets are used for evaluation: NW-UCLA [42],  
 261 NTU-RGBD60 (NTU60) [31], and NTU-RGBD120 (NTU120) [19]. The three datasets consist of  
 262 sequences with 1 or 2 subjects whose movement is measured in 3D space. For NW-UCLA, we  
 263 follow previous works [35] to train the model on view 1 and 2, and test the model on view 3. For  
 264 NTU60, we follow the official data split for the cross-subject (XSUB) and cross-view (XVIEW)  
 265 protocols. The similar is implemented on NTU120 for the cross-subject (XSUB) and cross-setting  
 266 (XSET) protocol. For ease of implementation, we project the 3D sequence into three orthonormal  
 267 2D spaces and use an MCAE defined on 2D space to process the three views of the sequences. Then  
 268 the segment activations from the three views are concatenated to form the representation. Four types

<sup>2</sup>Architectures of LSTM and 1D-Conv are detailed in the supplementary material.

Table 3: Performance (%) for unsupervised skeleton-based action classification. Column “Mod.” shows the data modality, where “S” indicates skeleton and “D” indicates depth map. Column “Cls.” shows the auxiliary classifier used for supervised training.

Model	Mod.	Cls.	NTU60		NTU120		NW-UCLA
			XSUB	XVIEW	XSUB	XSET	V1&V2 → V3
Luo <i>et al.</i> [21]	S+D	SLP	61.4	53.2	–	–	50.7
Li <i>et al.</i> [14]	S+D	SLP	68.1	63.9	–	–	62.5
SeBiReNet [23]	S	LSTM	–	79.7	–	–	80.3
LongT GAN [54]	S	SLP	39.1	48.1	–	–	74.3
MS <sup>2</sup> L [18]	S	SLP	52.6	–	–	–	76.8
CAE+ [27]	S	SLP	58.5	64.8	48.6	49.2	–
MCAE-MP (SLP)	S	SLP	<b>65.6</b>	74.7	<b>52.8</b>	<b>54.7</b>	83.6
P&C [35]	S	1-NN	50.7	76.1	–	–	<b>84.9</b>
MCAE-MP (1-NN)	S	1-NN	51.9	<b>82.4</b>	42.3	46.1	79.1

of disturbance are introduced for contrastive learning, namely jittering, spatial rotation, masking, and temporal smoothing. The readers are referred to the supplementary material for details.

The classification accuracy is put into three groups in Table 3. In the first group are the prior works that are not directly comparable as they use depth map [14, 21] or stronger auxiliary classifiers for supervised training [23]. In the second group, where our model is marked as MCAE-MP (SLP), a single layer perceptron (SLP) is trained as the auxiliary classifier with backbone parameters frozen. In the third group, where our model is marked as MCAE-MP (1NN), a 1-nearest-neighbor classifier is used instead of an SLP. Although MCAE-MP is a naive extension as it encodes joints separately and largely ignores their interactions, it achieves better or competitive performance compared with the baselines. Notably, on NTU60-XVIEW and NTU120-XSET where the training set and test set have different viewpoints, our model outperforms baselines by a clear margin thanks to the capsule-based representation which effectively captures viewpoint changes as transformations on input.

### 4.3 What does MCAE Learn?

To better understand what is encoded, we plot the learned snippet templates  $\mathcal{T}$  and segment templates  $\mathcal{P}$  in Fig. 4. Note that  $\mathcal{T}$  are initialized as random straight lines, and  $\mathcal{P}$  are initialized as arbitrary patterns composed randomly of  $\mathcal{T}$ . As shown in Fig. 4a, the snippets are mainly simple lines and hook-like curves that does not carry semantic information. Segment templates in Fig. 4b, however, bear some resemblance to the patterns shown in Fig. 3. This suggests that semantic-agnostic snippets are being aggregated into semantic-aware segments.

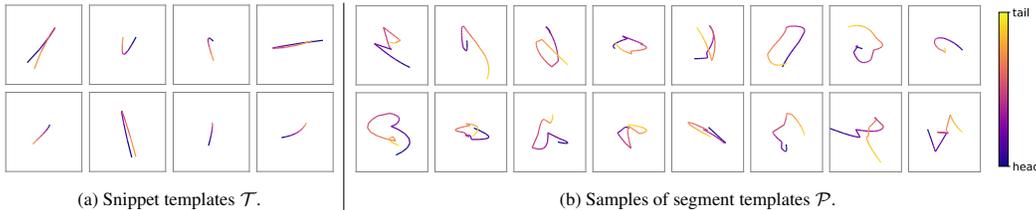


Figure 4: Templates learned from Trajectory20 dataset. Color indicates time.

We proceed to explore the information in SegCaps. In particular, we would like to see if SegCaps have learned transformation-invariant information. To this purpose, we randomly sample a trajectory from T20 dataset. The trajectory is first normalized so that its centroid is at (0, 0), then rotated clockwise by an angle  $\theta$ , and finally fed into the model. We examine the segment templates with the highest activation values (which reflects the trajectory’s semantics) and calculate the rotation angle  $\phi$  from those templates’ parameter  $\mathbf{B}$ . As shown in Table 4, the calculated  $\phi$  reveals two types of segments templates as we rotate the input. One type yields constant  $\phi$  (e.g. segment ID 2 for sample “absolute sine”), which indicates its rotation-invariance, the other has  $\phi$  that changes monotonically with  $\theta$  (e.g. segment ID 8 for sample “hexagon”), which shows its rotation-awareness. As for the activation values, samples from different categories activates different set of segments.

Table 4: Top-5 segment templates (sorted by segment activation  $\nu$  then segment ID for better visualization), and the rotation  $\phi$  calculated from their parameters  $B$ . Bold IDs are segments repeating across different  $\theta$ .

Input	$\theta = -10^\circ$		$\theta = -5^\circ$		$\theta = 0^\circ$		$\theta = 5^\circ$		$\theta = 10^\circ$	
	ID	$\phi$	ID	$\phi$	ID	$\phi$	ID	$\phi$	ID	$\phi$
	<b>2</b>	6.3	<b>2</b>	6.7	<b>2</b>	6.8	<b>2</b>	7.0	<b>2</b>	7.1
	<b>8</b>	6.9	<b>8</b>	9.0	<b>8</b>	11.2	<b>8</b>	13.9	<b>8</b>	16.5
	<b>12</b>	54.9	<b>12</b>	55.5	<b>12</b>	55.8	<b>12</b>	56.5	<b>12</b>	56.8
	<b>37</b>	-20.8	<b>37</b>	-19.8	<b>37</b>	-18.9	<b>37</b>	-17.9	<b>37</b>	-16.9
	<b>66</b>	50.2	<b>66</b>	52.5	<b>66</b>	55.4	<b>66</b>	59.0	<b>66</b>	62.4
	<b>2</b>	12.1	<b>2</b>	12.3	<b>2</b>	12.2	<b>2</b>	12.1	<b>2</b>	11.9
	<b>7</b>	8.2	<b>5</b>	-10.7	<b>5</b>	-10.1	<b>5</b>	-9.9	<b>7</b>	17.2
	<b>33</b>	65.1	<b>7</b>	10.7	<b>7</b>	13.4	<b>7</b>	15.4	<b>32</b>	-9.7
	<b>37</b>	-22.9	<b>37</b>	-22.3	<b>37</b>	-21.8	<b>37</b>	-21.3	<b>37</b>	-19.9
	<b>46</b>	45.7	<b>46</b>	47.5	<b>46</b>	48.6	<b>46</b>	50.2	<b>46</b>	51.6

Table 5: Top-5 segment templates (sorted by segment activation  $\nu$  then segment ID for better visualization), and the translation  $(x, y)$  calculated from their parameters  $B$ .

Input	$(\Delta x, \Delta y) = (-0.2, 0)$			$(\Delta x, \Delta y) = (-0.1, 0)$			$(\Delta x, \Delta y) = (0, 0)$			$(\Delta x, \Delta y) = (0, 0.1)$			$(\Delta x, \Delta y) = (0, 0.2)$		
	ID	$x$	$y$	ID	$x$	$y$	ID	$x$	$y$	ID	$x$	$y$	ID	$x$	$y$
	<b>2</b>	0.05	0.18	<b>2</b>	0.17	0.19	<b>2</b>	0.27	0.19	<b>2</b>	0.28	0.28	<b>2</b>	0.27	0.37
	<b>8</b>	0.01	-0.07	<b>8</b>	0.09	-0.06	<b>8</b>	0.18	-0.04	<b>8</b>	0.19	0.04	<b>8</b>	0.19	0.12
	<b>12</b>	-0.09	0.13	<b>12</b>	0.00	0.13	<b>12</b>	0.09	0.13	<b>12</b>	0.09	0.23	<b>12</b>	0.09	0.32
	<b>37</b>	0.10	-0.11	<b>37</b>	0.18	-0.11	<b>37</b>	0.27	-0.11	<b>37</b>	0.27	-0.03	<b>37</b>	0.27	0.05
	<b>66</b>	-0.12	0.16	<b>66</b>	-0.03	0.16	<b>66</b>	0.05	0.17	<b>66</b>	0.06	0.26	<b>66</b>	0.06	0.35
	<b>2</b>	0.04	0.2	<b>2</b>	0.14	0.19	<b>2</b>	0.24	0.19	<b>2</b>	0.24	0.28	<b>2</b>	0.23	0.38
	<b>5</b>	-0.01	0.30	<b>5</b>	0.07	0.29	<b>5</b>	0.16	0.29	<b>5</b>	0.16	0.38	<b>5</b>	0.15	0.46
	<b>7</b>	0.20	-0.16	<b>7</b>	0.28	-0.16	<b>7</b>	0.37	-0.15	<b>7</b>	0.36	-0.06	<b>7</b>	0.36	0.04
	<b>37</b>	0.04	-0.17	<b>37</b>	0.12	-0.16	<b>37</b>	0.21	-0.16	<b>37</b>	0.20	-0.07	<b>37</b>	0.20	0.01
	<b>46</b>	0.02	0.01	<b>46</b>	0.13	0.02	<b>46</b>	0.23	0.04	<b>46</b>	0.23	0.13	<b>46</b>	0.22	0.23

298 Meanwhile, the same sample under different rotation angle  $\theta$  gives stable segment activation, despite  
 299 some changes which are found to have no effect on the classification result.

300 We do a similar study on the translation component  $(x, y)$ , where we translate the input by  $(\Delta x, \Delta y)$ .  
 301 As shown in Table 5,  $(x, y)$  changes monotonically with  $(\Delta x, \Delta y)$  while the activated segments  
 302 remain stable. These results suggest that the semantics and transformation information has been  
 303 encoded separately in the segment activation  $\nu$  and transformation parameters  $B$ . In other words,  
 304 the encoded semantic information is robust against geometric transformations.

## 305 5 Conclusion

306 In this paper, we introduce MCAE, a framework that learns robust and discriminative representation  
 307 for keypoint motion. To resolve the intra-class variation of motion, we propose to learn a compact  
 308 and transformation-invariant motion representation using a two-level capsule-based representation  
 309 hierarchy. The efficacy of the learned representation is shown through an experimental study on  
 310 synthetic and real-world datasets. The output of MCAE could serve as mid-level representation  
 311 in other frameworks, e.g. Graph Convolution Network, for tasks that involve more context than  
 312 classification. We anticipate this work to inspire further studies that apply capsule-based models to  
 313 other time series processing tasks, such as joint modeling of visual appearance and motion in video.  
 314 The software and the T20 dataset of our research will be released to the research community.

315 Motion analysis techniques are in the foreground of the misuse of machine learning methods, among  
 316 which adverse societal impacts and privacy breach are two major concerns. Regarding the societal  
 317 impacts, admittedly, our method has both upside and downside. On one hand, a transformation-  
 318 invariant motion representation enables us better decode the information implicit in the trajectory,  
 319 which has applications for example in ethology. On the other hand, it could also be misused in mass  
 320 surveillance. Appropriate boundaries of use and ethical review are required to prevent potential  
 321 malicious applications. Regarding the privacy concerns, our method isolates the subjects' motion  
 322 from their sensitive information, such as gender and race.

## References

- [1] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):677–691, 2017.
- [2] Yong Du, Wei Wang, and Liang Wang. Hierarchical recurrent neural network for skeleton based action recognition. In *CVPR*, pages 1110–1118, 2015.
- [3] Kevin Duarte, Yogesh Rawat, and Mubarak Shah. VideoCapsuleNet: A simplified network for action detection. In *NeurIPS*, pages 7610–7619, 2018.
- [4] Tian Gan, Yongkang Wong, Daqing Zhang, and Mohan Kankanhalli. Temporal encoded f-formation system for social interaction detection. In *ACM Multimedia*, pages 937–946, 2013.
- [5] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social GAN: Socially acceptable trajectories with generative adversarial networks. In *CVPR*, pages 2255–2264, 2018.
- [6] Ahmed Refaat Hawas, Heba A. El-Khobby, Mohammed Abd-Elnaby, and Fathi E. Abd El-Samie. Gait identification by convolutional neural networks and optical flow. *Multimedia Tools and Applications*, 78(18):25873–25888, 2019.
- [7] Geoffrey E. Hinton, Alex Krizhevsky, and Sida D. Wang. Transforming auto-encoders. In *ICANN*, volume 6791 of *Lecture Notes in Computer Science*, pages 44–51. Springer, 2011.
- [8] Mariano Jaimez, Mohamed Souiai, Javier Gonzalez-Jimenez, and Daniel Cremers. A primal-dual framework for real-time dense RGB-D scene flow. In *ICRA*, pages 98–104, 2015.
- [9] Yanli Ji, Feixiang Xu, Yang Yang, Ning Xie, Heng Tao Shen, and Tatsuya Harada. Attention transfer (ANT) network for view-invariant action recognition. In *ACM Multimedia*, pages 574–582, 2019.
- [10] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [11] Adam Kosiorek, Sara Sabour, Yee Whye Teh, and Geoffrey E Hinton. Stacked capsule autoencoders. In *NeurIPS*, pages 15512–15522, 2019.
- [12] Mohamed Ilyes Lakkhal, Oswald Lanz, and Andrea Cavallaro. View-1stm: Novel-view video synthesis through view decomposition. In *ICCV*, pages 7576–7586, 2019.
- [13] Binlong Li, Octavia I. Camps, and Mario Sznaiar. Cross-view activity recognition using hanklets. In *CVPR*, pages 1362–1369, 2012.
- [14] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S. Kankanhalli. Unsupervised learning of view-invariant action representations. In *NeurIPS*, pages 1262–1272, 2018.
- [15] Maosen Li, Siheng Chen, Xu Chen, Ya Zhang, Yanfeng Wang, and Qi Tian. Actional-structural graph convolutional networks for skeleton-based action recognition. In *CVPR*, pages 3595–3603, 2019.
- [16] Wanqing Li, Zhengyou Zhang, and Zicheng Liu. Action recognition based on a bag of 3D points. In *CVPR Workshops*, pages 9–14, 2010.
- [17] Yan Li, Bin Ji, Xintian Shi, Jianguo Zhang, Bin Kang, and Limin Wang. TEA: Temporal excitation and aggregation for action recognition. In *CVPR*, pages 906–915, 2020.
- [18] Lilang Lin, Sijie Song, Wenhan Yang, and Jiaying Liu. MS2L: multi-task self-supervised learning for skeleton based action recognition. In *ACM Multimedia*, pages 2490–2498, 2020.
- [19] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C. Kot. NTU RGB+D 120: A large-scale benchmark for 3D human activity understanding. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 42(10):2684–2701, 2020.

- 369 [20] Ziyu Liu, Hongwen Zhang, Zhenghao Chen, Zhiyong Wang, and Wanli Ouyang. Disentangling  
370 and unifying graph convolutions for skeleton-based action recognition. In *CVPR*, pages 143–  
371 152, 2020.
- 372 [21] Zelun Luo, Boya Peng, De-An Huang, Alexandre Alahi, and Li Fei-Fei. Unsupervised learn-  
373 ing of long-term motion dynamics for videos. In *CVPR*, pages 7101–7110. IEEE Computer  
374 Society, 2017.
- 375 [22] Joe Yue-Hei Ng, Matthew J. Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat  
376 Monga, and George Toderici. Beyond short snippets: Deep networks for video classification.  
377 In *CVPR*, pages 4694–4702, 2015.
- 378 [23] Qiang Nie, Ziwei Liu, and Yunhui Liu. Unsupervised 3d human pose representation with  
379 viewpoint and pose disentanglement. In *ECCV*, volume 12364 of *Lecture Notes in Computer  
380 Science*, pages 102–118. Springer, 2020.
- 381 [24] Dan Oneata, Jakob Verbeek, and Cordelia Schmid. Action and event recognition with fisher  
382 vectors on a compact feature set. In *ICCV*, pages 1817–1824, 2013.
- 383 [25] John Paparrizos and Luis Gravano. k-shape: Efficient and accurate clustering of time series.  
384 In *SIGMOD*, pages 1855–1870, 2015.
- 385 [26] Vasu Parameswaran and Rama Chellappa. View invariance for human action recognition.  
386 *International Journal of Computer Vision*, 66(1):83–101, 2006.
- 387 [27] Haocong Rao, Shihao Xu, Xiping Hu, Jun Cheng, and Bin Hu. Augmented skeleton based  
388 contrastive action learning with momentum LSTM for unsupervised action recognition. *Infor-  
389 mation Sciences*, 569:90–109, August 2021.
- 390 [28] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In  
391 *NIPS*, pages 3856–3866, 2017.
- 392 [29] Hiroaki Sakoe and Seibi Chiba. Dynamic programming algorithm optimization for spoken  
393 word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26(1):43–  
394 49, 1978.
- 395 [30] Arun Sankisa, Arjun Punjabi, and Aggelos K. Katsaggelos. Temporal capsule networks for  
396 video motion estimation and error concealment. *Signal Image Video Process.*, 14(7):1369–  
397 1377, 2020.
- 398 [31] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. NTU RGB+D: A large scale  
399 dataset for 3D human activity analysis. In *CVPR*, pages 1010–1019, 2016.
- 400 [32] Xingjian Shi, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun  
401 Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcast-  
402 ing. In *NIPS*, pages 802–810, 2015.
- 403 [33] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recog-  
404 nition in videos. In *NIPS*, pages 568–576, 2014.
- 405 [34] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of  
406 video representations using LSTMs. In *ICML*, volume 37 of *JMLR Workshop and Conference  
407 Proceedings*, pages 843–852, 2015.
- 408 [35] Kun Su, Xiulong Liu, and Eli Shlizerman. PREDICT & CLUSTER: unsupervised skeleton  
409 based action recognition. In *CVPR*, pages 9628–9637, 2020.
- 410 [36] Romain Tavenard, Johann Faouzi, Gilles Vandewiele, Felix Divo, Guillaume Androz, Chester  
411 Holtz, Marie Payne, Roman Yurchak, Marc Rußwurm, Kushal Kolar, and Eli Woods. Tslern,  
412 a machine learning toolkit for time series data. *Journal of Machine Learning Research*,  
413 21(118):1–6, 2020.
- 414 [37] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning  
415 spatiotemporal features with 3D convolutional networks. In *ICCV*, pages 4489–4497, 2015.

- 416 [38] Raviteja Vemulapalli, Felipe Arrate, and Rama Chellappa. Human action recognition by rep-  
417 resenting 3D skeletons as points in a Lie group. In *CVPR*, pages 588–595, 2014.
- 418 [39] Shruti Vyas, Yogesh Singh Rawat, and Mubarak Shah. Multi-view action recognition using  
419 cross-view video prediction. In *ECCV*, volume 12372 of *Lecture Notes in Computer Science*,  
420 pages 427–444. Springer, 2020.
- 421 [40] Heng Wang and Cordelia Schmid. Action recognition with improved trajectories. In *ICCV*,  
422 pages 3551–3558, 2013.
- 423 [41] Jiang Wang, Zicheng Liu, Ying Wu, and Junsong Yuan. Mining actionlet ensemble for action  
424 recognition with depth cameras. In *CVPR*, pages 1290–1297, 2012.
- 425 [42] Jiang Wang, Xiaohan Nie, Yin Xia, Ying Wu, and Song-Chun Zhu. Cross-view action model-  
426 ing, learning, and recognition. In *CVPR*, pages 2649–2656, 2014.
- 427 [43] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks.  
428 In *CVPR*, pages 7794–7803, 2018.
- 429 [44] Yunbo Wang, Lu Jiang, Ming-Hsuan Yang, Li-Jia Li, Mingsheng Long, and Li Fei-Fei. Eidetic  
430 3D LSTM: A model for video prediction and beyond. In *ICLR*, 2019.
- 431 [45] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for  
432 skeleton-based action recognition. In *AAAI*, 2018.
- 433 [46] Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network  
434 for action recognition. In *CVPR*, pages 588–597, 2020.
- 435 [47] Yue Yu, Niehao Tian, Xiangru Chen, and Ying Li. Skeleton capsule net: An efficient network  
436 for action recognition. In *ICVRV*, pages 74–77, 2018.
- 437 [48] Pengfei Zhang, Cuiling Lan, Wenjun Zeng, Junliang Xing, Jianru Xue, and Nanning Zheng.  
438 Semantics-guided neural networks for efficient skeleton-based human action recognition. In  
439 *CVPR*, pages 1112–1121, 2020.
- 440 [49] Xikun Zhang, Chang Xu, and Dacheng Tao. Context aware graph convolution for skeleton-  
441 based action recognition. In *CVPR*, pages 14333–14342, 2020.
- 442 [50] Yiyi Zhang, Li Niu, Ziqi Pan, Meichao Luo, Jianfu Zhang, Dawei Cheng, and Liqing Zhang.  
443 Exploiting motion information from unlabeled videos for static image action recognition. In  
444 *AAAI*, pages 12918–12925, 2020.
- 445 [51] Yujia Zhang, Xiaodan Liang, Dingwen Zhang, Min Tan, and Eric P Xing. Unsupervised  
446 object-level video summarization with online motion auto-encoder. *Pattern Recognition Let-  
447 ters*, 130:376–385, 2020.
- 448 [52] Yongheng Zhao, Tolga Birdal, Haowen Deng, and Federico Tombari. 3D point capsule net-  
449 works. In *CVPR*, pages 1009–1018, 2019.
- 450 [53] Yongheng Zhao, Tolga Birdal, Jan Eric Lenssen, Emanuele Menegatti, Leonidas J. Guibas, and  
451 Federico Tombari. Quaternion equivariant capsule networks for 3D point clouds. In *ECCV*,  
452 volume 12346 of *Lecture Notes in Computer Science*, pages 1–19. Springer, 2020.
- 453 [54] Nenggan Zheng, Jun Wen, Risheng Liu, Liangqu Long, Jianhua Dai, and Zhefeng Gong. Un-  
454 supervised representation learning with long-term dynamics for skeleton based action recog-  
455 nition. In *AAAI*, pages 2644–2651, 2018.
- 456 [55] Tao Zhuo, Zhiyong Cheng, Peng Zhang, Yongkang Wong, and Mohan Kankanhalli. Explain-  
457 able video action reasoning via prior knowledge and state transitions. In *ACM Multimedia*,  
458 pages 521–529, 2019.

459 **Checklist**

- 460 1. For all authors...
- 461 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
462 contributions and scope? [Yes] Please see Section 1.
- 463 (b) Did you describe the limitations of your work? [Yes]
- 464 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Please  
465 check Section 5.
- 466 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
467 them? [Yes]
- 468 2. If you are including theoretical results...
- 469 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 470 (b) Did you include complete proofs of all theoretical results? [N/A]
- 471 3. If you ran experiments...
- 472 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
473 imental results (either in the supplemental material or as a URL)? [No] The source  
474 code and dataset will be released with the publication of the paper. Implementation  
475 details are covered in the supplementary material.
- 476 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
477 were chosen)? [Yes] Training details are mostly provided in Section 4 while some  
478 extra information is in the supplementary material.
- 479 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
480 ments multiple times)? [Yes] We provided error bar on our T20 dataset. For experi-  
481 ments on NW-UCLA, NTURGBD-60, and NTURGBD-120, we followed the protocol  
482 in prior works.
- 483 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
484 of GPUs, internal cluster, or cloud provider)? [Yes] Please see the first paragraph of  
485 Section 4.
- 486 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 487 (a) If your work uses existing assets, did you cite the creators? [Yes] We used the publicly  
488 available NTURGBD60/120 and NW-UCLA dataset. The original papers have been  
489 cited.
- 490 (b) Did you mention the license of the assets? [N/A]
- 491 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
492 The T20 dataset will be released with the publication of the paper.
- 493 (d) Did you discuss whether and how consent was obtained from people whose data  
494 you're using/curating? [Yes] The data is publicly available and the consent has been  
495 described in the published paper.
- 496 (e) Did you discuss whether the data you are using/curating contains personally identifi-  
497 able information or offensive content? [N/A] No such information was found in the  
498 dataset we use.
- 499 5. If you used crowdsourcing or conducted research with human subjects...
- 500 (a) Did you include the full text of instructions given to participants and screenshots, if  
501 applicable? [N/A] No crowdsourcing research was conducted.
- 502 (b) Did you describe any potential participant risks, with links to Institutional Review  
503 Board (IRB) approvals, if applicable? [N/A]
- 504 (c) Did you include the estimated hourly wage paid to participants and the total amount  
505 spent on participant compensation? [N/A]