

A Novel Hierarchical Framework for Human Action Recognition

Hongzhao Chen^a, Guijin Wang^{a,*}, Jing-Hao Xue^b, Li He^a

^a*Department of Electronic Engineering, Tsinghua University, Beijing 100084, China*

^b*Department of Statistical Science, University College London, London WC1E 6BT, UK*

Abstract

In this paper, we propose a novel two-level hierarchical framework for three-dimensional (3D) skeleton-based action recognition, in order to tackle the challenges of high intra-class variance, movement speed variability and high computational costs of action recognition. In the first level, a new part-based clustering module is proposed. In this module, we introduce a part-based five-dimensional (5D) feature vector to explore the most relevant joints of body parts in each action sequence, upon which action sequences are automatically clustered and the high intra-class variance is mitigated. In the second level, there are two modules, motion feature extraction and action graphs. In the module of motion feature extraction, we utilize the cluster-relevant joints only and present a new statistical principle to decide the time scale of motion features, to reduce computational costs and adapt to variable movement speed. In the action graphs module, we exploit these 3D skeleton-based motion features to build action graphs, and devise a new score function

*Corresponding author. Tel.: +86-18911389502; Fax: +86-62770317

Email addresses: jordanchan1004@163.com (Hongzhao Chen), wangguijin@tsinghua.edu.cn (Guijin Wang), jinghao.xue@ucl.ac.uk (Jing-Hao Xue), hhappy06@gmail.com (Li He)

based on maximum-likelihood estimation for action graph-based recognition. Experiments on the Microsoft Research Action3D dataset and the University of Texas Kinect Action dataset demonstrate that our method is superior or at least comparable to other state-of-the-art methods, achieving 95.56% recognition rate on the former dataset and 95.96% on the latter one.

Keywords: action recognition, 3D skeleton, hierarchical framework, part-based, time scale, action graphs.

1. Introduction

Action recognition is an active research topic that focuses on labeling a motion sequence as one of the known actions. It can be widely applied in human-computer interaction, health care, video surveillance, etc. In order to achieve high accuracy and great robustness for real-world applications, an action recognition system has to overcome three challenges: high intra-class variance with low inter-class variance, variable movement speed, and high computational costs. As shown in Fig. 1, people may perform the same action of *Side Boxing* in quite different ways, by using one hand or two hands, leading to high intra-class variance. Meanwhile, people may also perform the same action with variable movement speed, as demonstrated in Fig. 2.

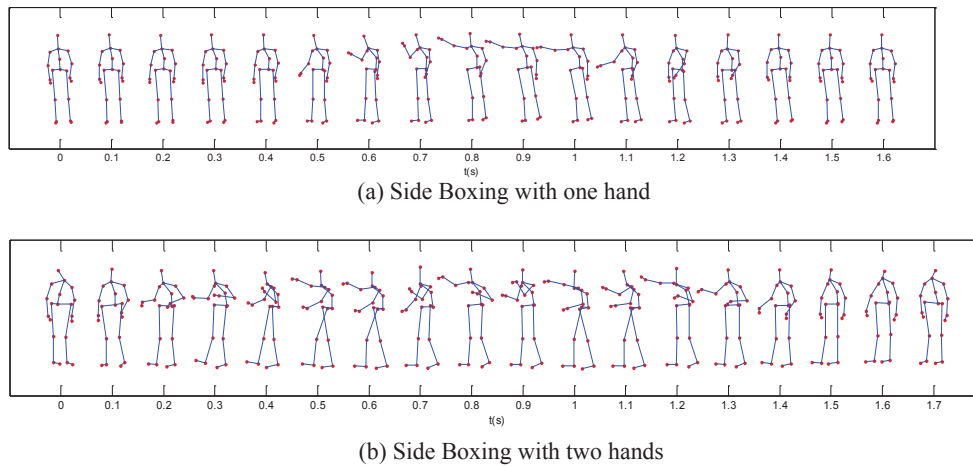


Figure 1: An illustrative example of high intra-class variance. Two panels present skeleton sequence diagrams of action *Side Boxing* sampled at 10fps.

Prior to 2010, many color image-based methods of action classification

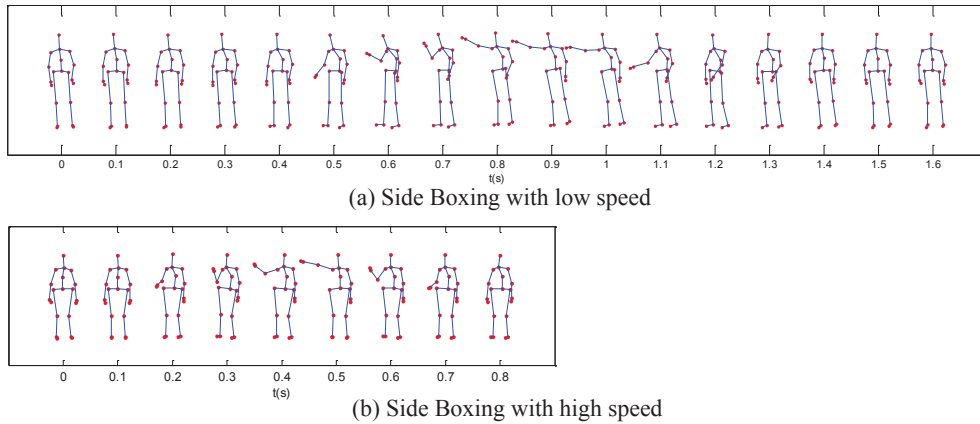


Figure 2: An illustrative examples of variable movement speed. Two panels present skeleton sequence diagrams of action *Side Boxing* sampled at 10fps.

13 had been studied [1]. However, these methods have relatively low recognition
 14 accuracies and thus they are unable to be applied in real-world applications.

15 The situation has been much improved as technologies on depth imag-
 16 ing advance quickly [2–5]. Recent works of action recognition could be di-
 17 vided into two types, depth map-based methods [6–10] and 3D skeleton-based
 18 methods [11–22]. The former directly takes sequences of depth maps as in-
 19 put, while the latter utilizes 3D skeleton sequences inferred from depth maps.
 20 Fig. 3 shows the color image, depth image and 3D skeleton acquired from a
 21 Kinect sensor.

22 Depth map-based methods extract features from depth maps to describe
 23 the human poses and model the transition of poses. The widely-used features
 24 include sampled 3D points from silhouettes [6, 7], histograms of oriented
 25 gradients [8], histogram of oriented 4D normals [9], histogram of oriented
 26 principal components [10], etc. However, the extraction of these features is

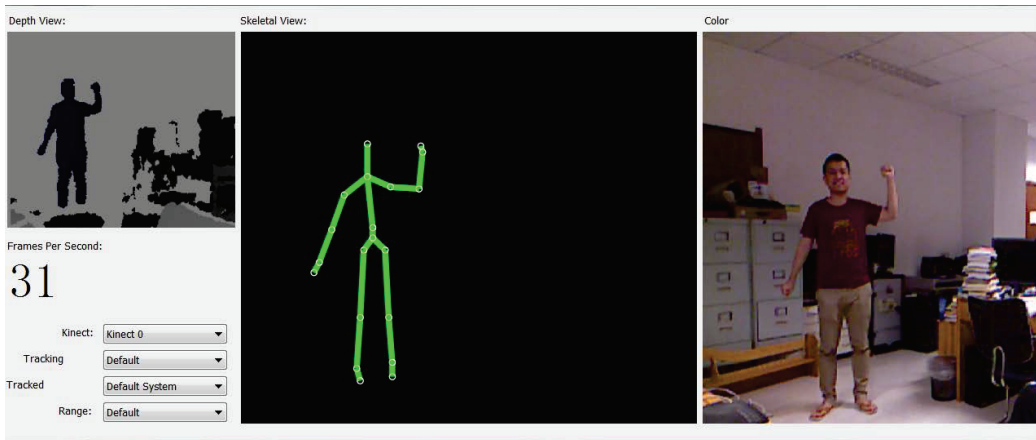


Figure 3: Depth image, skeleton and color image.

27 often time-consuming, making them hardly applicable in real-time scenarios.

28 In fact 3D human skeleton, which could be reliably estimated from depth
 29 maps in real time [23–25], is an efficient and concise surrogate to describe the
 30 human poses. In most 3D skeleton-based methods [11, 12, 14, 16, 18, 21], mo-
 31 tion features are represented by pair-wise differences of joint positions within
 32 the current frame or between the current frame and the previous frames.
 33 Hence motion features extracted from 3D skeleton can efficiently model the
 34 action dynamics. Kapsouras *et al.* [21] further considered the time scale for
 35 motion features to fit various movement speeds. However, no principle has
 36 been supplied yet for how to determine the time scale. Some other histogram-
 37 based features are also proposed, like histograms of 3D joints [13], space time
 38 pose [17], histogram of oriented displacements [15], points in a lie group [19],
 39 etc.

40 Yang *et al.* [12] used a 3D skeleton-based method to achieve higher recog-
 41 nition rates than the depth map-based method [6], but in their method it was

42 still difficult to distinguish similar actions, such as *Draw X*, *Draw Tick* and
43 *Draw Circle*. Mining techniques were adopted for finding relevant joints for
44 each action [11, 14], which indeed improved the classification performance.
45 However, mining techniques are computationally expensive and difficult to
46 be expanded to new actions. Chen *et al.* [18] proposed a simple and effec-
47 tive hierarchical model to cluster similar actions and improved the accuracy
48 of recognition. However, the hierarchical model therein was manually con-
49 structed, which limited the expansibility of the method. A new representa-
50 tion called SMIJ (sequence of the most informative joints) was proposed to
51 select the most informative joints for performing an action [22], which was
52 easy to interpret. However, this representation required segmenting each ac-
53 tion sequence into several windows beforehand, which raised a difficult prob-
54 lem of choosing the temporal window size. Moreover, it performed poorly
55 when the skeleton data were noisy or the actions were based on almost the
56 same joints.

57 In this paper, in order to tackle the three challenges aforementioned, we
58 propose a novel two-level hierarchical framework for 3D skeleton-based action
59 recognition. The first level of the framework consists of a part-based cluster-
60 ing module. In this module, a part-based five-dimensional feature vector is
61 introduced to explore the most relevant joints of body parts in each action
62 sequence, upon which action sequences are clustered. Distinct sequences of
63 the same action could be grouped into different clusters, enabling us to cope
64 with the problem of high intra-class variance. Hence this module groups
65 similar actions together and divides the recognition task for various actions
66 into several smaller and simpler tasks, which can significantly improve the

67 final performance (for example by more than 10% in our first experiment).
68 In the second level of the framework, there are two modules, motion feature
69 extraction and action graphs. For each cluster, only the relevant joints ob-
70 tained from the first level are utilized for motion feature extraction, which
71 not only enhances the validity of the extracted features but also reduces the
72 computational costs remarkably. We also investigate and derive a statistical
73 principle for determining the time scale of motion features to deal with the
74 problem of variable movement speed. After motion feature extraction, we
75 apply action graphs to these 3D skeleton-based motion features to finally
76 classify actions. Experiments on the Microsoft Research Action3D dataset
77 and the UTKinect-Action dataset show that our method is noticeably su-
78 perior or at least comparable to other state-of-the-art methods, achieving
79 recognition rates of 95.56% and 95.96%, respectively on the two datasets.

80 The remainder of this paper is organized as follows. In Section 2 the
81 proposed hierarchical framework is described. Details of the three key mod-
82 ules, part-based clustering, motion feature extraction and action graphs, are
83 followed in Section 3. Experimental performance of our method is demon-
84 strated in Section 4. Section 5 concludes our work and discusses the future
85 work.

86 **2. Hierarchical Framework**

87 As shown in Fig. 4, our hierarchical framework for 3D skeleton-based
88 action recognition consists of three modules: part-based clustering, motion
89 feature extraction and action graphs classification. For an action sequence,
90 the hierarchical framework first decides its cluster. Motion features are then

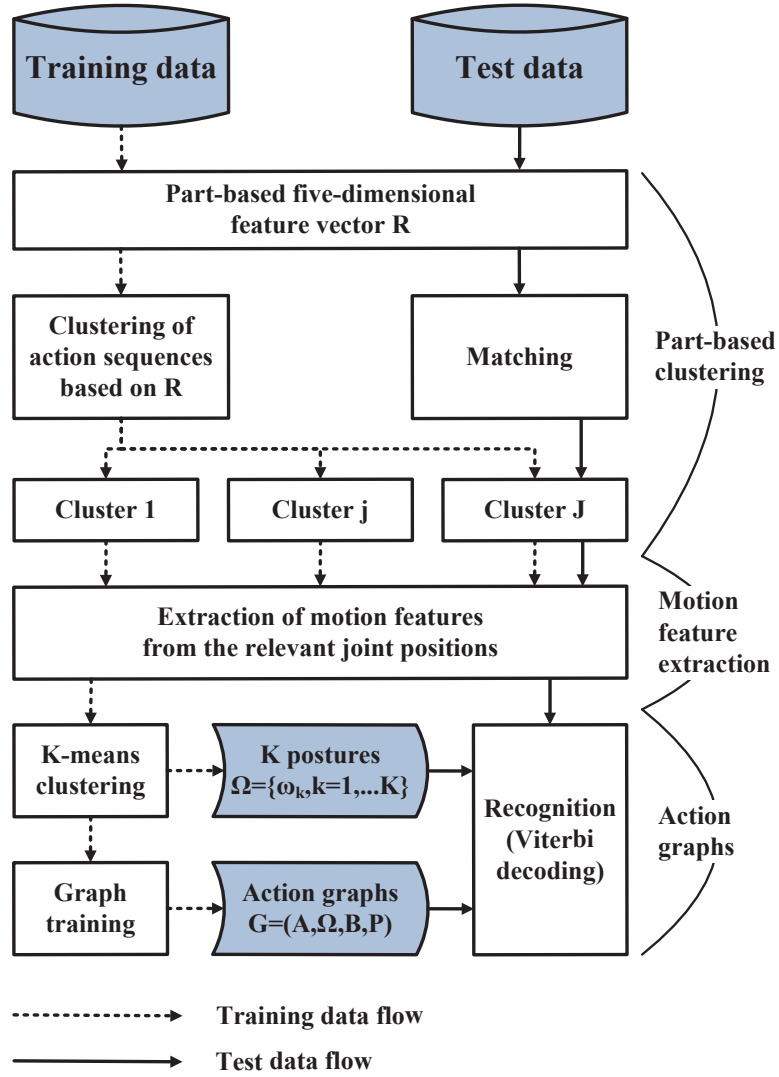


Figure 4: The hierarchical framework of our proposed method.

91 extracted from the relative joints. In the final classification, we apply the
92 Viterbi decoding to the action graphs. Unlike [18], the hierarchical framework
93 here is automatically learned from data during the training procedure with-
94 out manual intervention. Correspondingly, our framework has three main
95 characteristics.

96 Firstly, a part-based five-dimensional feature vector is introduced to ex-
97 plore the most relevant joints of body parts in each action sequence. Then
98 action sequences are clustered to construct the first level of the hierarchical
99 model. More details will be described in Section 3.1. The part-based clus-
100 tering module groups similar actions together and partitions the recognition
101 task for various actions into several smaller and simpler tasks, which could
102 significantly improve the final performance, as verified in Section 4.3.1.

103 Secondly, for each cluster, only the relevant joints are utilized for mo-
104 tion feature extraction. Since irrelevant joints provide little information for
105 classification, motion features of the relevant joints can not only enhance the
106 validity of the extracted features but also reduce the computational costs. We
107 extract motion features with higher order and time scale (see Section 3.2),
108 which improves feature representability and thus achieves higher accuracy
109 as demonstrated in Section 4.3.2. We further investigate how to choose the
110 time scale for any specified dataset.

111 Thirdly, we apply action graphs for classification. Unlike [6], our action
112 graphs are applied to skeleton-based motion features. Postures are obtained
113 via K-means clustering, upon which action graphs are trained. From max-
114 imum likelihood estimation, we derive a new score function for the Viterbi
115 decoding. We also investigate action graphs for early detection of actions,

116 which will be discussed in Section 4.3.3.

117 **3. Modules**

118 *3.1. Part-based Clustering*

119 The joints of a human body could be divided into several parts, as actions
120 are only related to certain parts of the body. Given the relevance of parts
121 for each action sequence, actions that are relevant to different joints could be
122 discriminated easily. So we could cluster action sequences with same relevant
123 parts together and divide the recognition task into several simpler tasks in a
124 cluster. In our method, a part-based five-dimensional feature vector is defined
125 to explore the most relevant joints to body parts in each action sequence,
126 then action sequences are clustered by using these features to construct the
127 first level of the hierarchical framework.

128 We define a body part as a set of joints close to each other. For clarity and
129 simplicity, here we assume that the number of joints involved is 20 according
130 to the Kinect sensor. Other situations with different numbers of joints could
131 be adapted without difficulty. Here five body parts and the joints relevant
132 to each body part are defined as

$$\begin{aligned}
O_1 = LUE &= \{left\ shoulder(1), left\ elbow(2), \\
&\quad left\ wrist(3), left\ hand(4)\}; \\
O_2 = RUE &= \{right\ shoulder(5), right\ elbow(6), \\
&\quad right\ wrist(7), right\ hand(8)\}; \\
O_3 = LLE &= \{left\ hip(9), left\ knee(10), \\
&\quad left\ ankle(11), left\ foot(12)\}; \\
O_4 = RLE &= \{right\ hip(13), right\ knee(14), \\
&\quad right\ ankle(15), right\ foot(16)\}; \\
O_5 = TRS &= \{head(17), shoulder\ center(18), \\
&\quad spine(19), hip\ center(20)\}.
\end{aligned}$$

133 To measure the relevance of the five body parts to an action sequence,
134 we construct a part-based five-dimensional feature vector $R = [R_1, \dots, R_5]$.
135 Given a T -frame sequence of joints positions $X = \{X^1, \dots, X^T\}$, where
136 $X^t = [x_1^t, \dots, x_{20}^t]$ in which x_i^t is the 3D coordinate of the i th joint in the t th
137 frame, the variance vector $V = [V_1, \dots, V_{20}]$ of the joints in the sequence is
138 calculated as

$$V_i = \sum_{t=1}^T \|x_i^t - \bar{x}_i\|^2,$$

139 where the mean coordinate $\bar{x}_i = \frac{1}{T} \sum_{t=1}^T x_i^t$. For incomplete skeletons, the
140 corresponding entries V_i can be simply set as zero for the missing joints.
141 Then the feature vector R , which represents the relevance of body parts to

142 an action sequence, is defined as:

$$R_j = 1_{\{(\hat{R}_j / \max_{l \in \{1, \dots, 5\}} \hat{R}_l) > \eta\}}, \text{ where } \hat{R}_j = \sum_{i \in O_j} V_i. \quad (1)$$

143 Here $1_{\{\cdot\}}$ is the indicator function valued in $\{0, 1\}$; the threshold η can be
 144 decided by cross-validation using the training set.

145 Then action sequences can be clustered based on R , which automatically
 146 constructs the first level of the hierarchical framework during the training
 147 phase. A practical example of part-based clustering is shown in Fig. 5. As
 148 different people may perform the same action in different ways (actions with
 149 large intra-class variances), distinct sequences of the same action are allowed
 150 to be grouped into different clusters. For example, one person may perform
 151 action *Side Boxing* with two hands while another person may do it with only
 152 one hand, which means that R for these two sequences are different so that
 153 they may belong to different clusters. Allowing distinct sequences of the
 154 same action to be grouped into different clusters improves the performance,
 155 as verified in Section 4.3.1.

156 Our method is also readily expansible for adding new actions. To add a
 157 new action to the dataset, we could simply calculate R for sequences of this
 158 new action, then join them to proper clusters or create a new cluster, and
 159 finally retrain the action graphs for the actions of the affected clusters.

160 3.2. Motion Feature Extraction

161 For each cluster, we extract motion features F^t , $t = 1, \dots, T$, from the
 162 relevant joints $\hat{X}^t = \{x_i^t \mid i \in O_j, R_j = 1, j = 1, \dots, 5\}$ only. For exam-
 163 ple, assume $R = [1, 1, 0, 0, 0]$ for a cluster, which means that the actions in
 164 this cluster are mainly relative to the two hands. Then the relevant joints

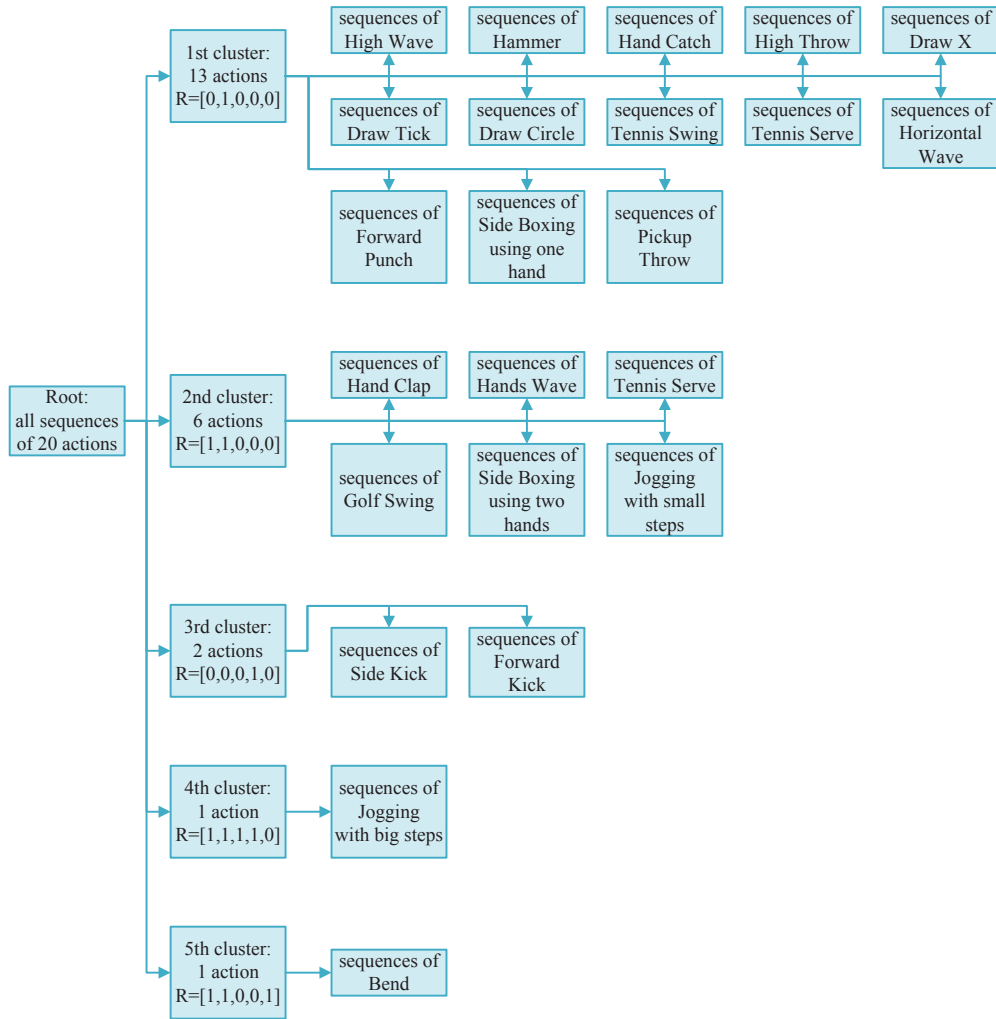


Figure 5: The part-based clustering result obtained from the 20 actions.

165 \hat{X}^t are the joints of $O_1(LUE)$ and $O_2(RUE)$, i.e., $\hat{X}^t = \{x_1^t, x_2^t, \dots, x_8^t\}$.
 166 As we extract motion features only from relevant joints, the dimension of
 167 \hat{X}^t is reduced. Since joints of irrelevant parts provide little information for
 168 identification of actions, the motion features extracted from only the joints
 169 of relevant parts can not only enhance the validity of the features but also
 170 greatly reduce the computational costs. The procedure of motion feature
 171 extraction is described as follow.

172 Similarly to previous works [12, 21], we compute pair-wise differences
 173 of joint positions between specific frames to obtain motion features. These
 174 features consist of three components: static pose SP , dynamic pose DP and
 175 offset pose OP . The static pose represents the static state at the current
 176 frame, the dynamic pose represents the instant motion, and the offset pose
 177 represents the offset from the initial state.

178 In [12], the higher order terms (position differences of different joints)
 179 were utilized to reduce noise but without considering the time scale. In
 180 contrast, [21] considered the time scale but ignored the higher order terms.
 181 In this paper, we combine them together, which means: We calculate SP ,
 182 OP and DP with higher order terms as in [12], while for DP , to take the
 183 time scale into consideration, we calculate it with three previous frames (1,
 184 5, 10 frames) before the current frame, respectively. This combination aims
 185 to capture the dynamics in various time scales and to reduce noise at the
 186 same time, which could enhance representability of the feature. We further
 187 investigate and devise a statistical principle for choosing the time scale, i.e.,
 188 the number of previous frames used. It will be demonstrated that, given
 189 a dataset, using the time scale up to the standard deviations of sequence

190 lengths is highly possible to offer an optimal performance. More details
 191 could be found in Section 4.3.2.

192 To extract the motion features for the t th frame of a T -frame sequence,
 193 we calculate the three components from only the relevant joint positions
 194 $\hat{X}^t = \{\hat{x}_m^t \mid m = 1, \dots, M\}$, where M is the number of relevant joints:

$$SP^t = \{\hat{x}_i^t - \hat{x}_j^t \mid i, j = 1, \dots, M; i \neq j\} ,$$

$$DP^t = \{\hat{x}_i^t - \hat{x}_j^{t-s} \mid i, j = 1, \dots, M; s = 1, 5, 10\} ,$$

$$OP^t = \{\hat{x}_i^t - \hat{x}_j^1 \mid i, j = 1, \dots, M\} .$$

195 Then we concatenate the three components to obtain the motion features
 196 $F_{ori}^t = [SP^t, DP^t, OP^t]$. After that, we apply principle component analysis
 197 (PCA) to reduce the dimension of the features and reach the final motion
 198 features

$$F^t = W_{opt}(F_{ori}^t - \mu) ,$$

199 where W_{opt} is the optimal projection matrix and μ is the mean of all F_{ori}^t .
 200 Note that here PCA is utilized to reduce computational costs and the result-
 201 ing motion features are later fed into action graphs for classification rather
 202 than directly used for classification. Thus, supervised dimension-reduction
 203 techniques like linear discriminant analysis (LDA) are not suitable here. De-
 204 note the dimension of the final motion features F^t as L . The impact of L on
 205 recognition performance will be investigated by experiments in Section 4.4.

206 3.3. Action Graphs

207 Many classification methods have been introduced to action recognition,
208 among which action graphs [26] are an effective method to explicitly model
209 the action dynamics. Action graphs were applied to depth map-based action
210 recognition in [6]. Unlike [6], here we apply action graphs to 3D skeleton-
211 based features and derive a new score function for recognition.

212 An action can be represented by transitions of several postures, here a
213 posture means the similar motion features for the frames in a salient state.
214 Thus we could model an action as a weighted directed graph, whose nodes
215 represent the postures of the action and whose edges represent the transi-
216 tional probabilities between two postures.

217 Fig. 6 is a sketch of two action graphs. Action *High Throw* consists of
218 three postures, $\omega_1 = \textit{hand lifting}$, $\omega_2 = \textit{throwing}$ and $\omega_3 = \textit{hand putting down}$.
219 The notation $p_n[i, j]$ near the arrows represent the transitional probabilities
220 between postures ω_i and ω_j . In most cases of action *High Throw*, it starts
221 from *hand lifting*, transits to *throwing* and ends at *hand putting down*. It
222 is also quite possible that two consecutive frames stay at the same posture,
223 which is represented as a self-loop edge in the action graph. Situations are
224 similar for action *High Wave*, which shares two postures ω_1 and ω_3 with
225 the former action and has another two new postures $\omega_4 = \textit{wave slightly}$ and
226 $\omega_5 = \textit{wave substantially}$. Nevertheless, because of high intra-class variance,
227 there are two or more possible paths (e.g. 1-4-3 or 1-5-3) for the same action.

228 As actions often share postures, we could obtain postures by clustering
229 the motion features of all frames from all training samples via an algorithm
230 like K-means clustering, with each cluster center representing a posture.

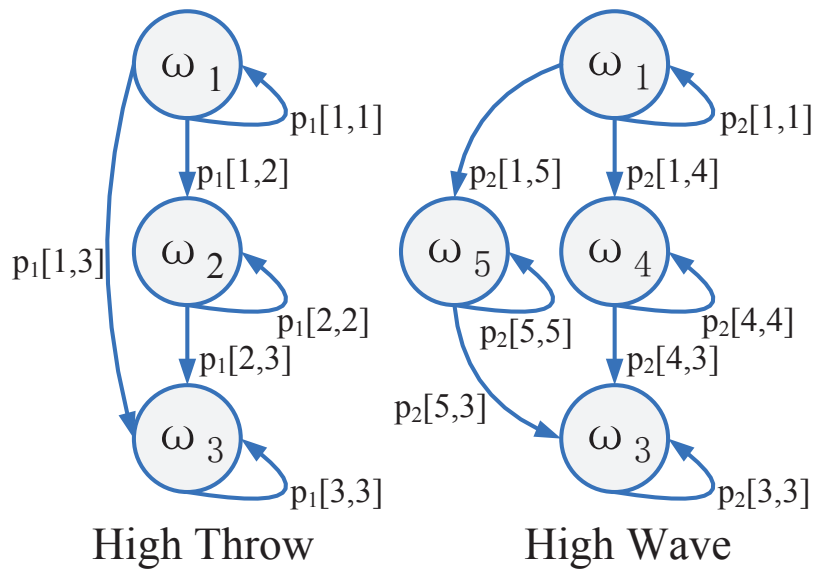


Figure 6: A sketch of two action graphs. In the sketch, each ω_i represents a posture and the notation $p_n[i, j]$ near the arrows represent the transitional probabilities between postures ω_i and ω_j for the n th action. This sketch intends to visualize the relationship of postures, actions and transitional probabilities.

231 Consider an action set of N actions $A = \{A_n\}_{n=1,\dots,N}$ that contains K
 232 postures $\Omega = \{\omega_k\}_{k=1,\dots,K}$, we could model them as a set of weighted directed
 233 graphs, which can be represented by a quadruplet $G = (A, \Omega, B, P)$, where

$$A = \{A_n\}_{n=1,\dots,N} ,$$

$$\Omega = \{\omega_k\}_{k=1,\dots,K} ,$$

$$B = \{B_n\}_{n=1,\dots,N} ,$$

$$P = \{P_k\}_{k=1,\dots,K} ,$$

234 in which $B_n = \{p(\omega_j|\omega_i, A_n)\}_{i,j=1,\dots,K}$, for $n = 1, \dots, N$, is the transitional
 235 probability matrix of the n th action A_n , and $P_k(F^t) = p(F^t|\omega_k)$ is the con-
 236 ditional probability of an observation F^t to be generated from the node
 237 ω_k . We assume that the distribution of the observations for a node can
 238 be approximated by an isotropic normal distribution: $P_k(F^t) = p(F^t|\omega_k) =$
 239 $\frac{1}{(2\pi)^{L/2}\sigma^L} \exp(-\frac{1}{2\sigma^2}\|F^t - \omega_k\|^2)$. Here σ is the standard deviation for all L
 240 dimensions. Such action graphs G can be trained as described in [26].

241 For classification, we could apply maximum likelihood estimation. Given
 242 an action sequence of T frames, let F^t be the final motion feature vector
 243 for the t th frame, and $F = \{F^t\}_{t=1,\dots,T}$ be the motion feature sequence.
 244 The posture sequence corresponding to F is denoted by $S = \{S^t\}_{t=1,\dots,T}$,
 245 where $S^t \in \Omega$ for all t . The recognition of the most likely action A^* that
 246 generates the observation F can be then formulated as a maximum likelihood

247 estimation:

$$\begin{aligned}
A^* &= \arg \max_{A_n \in A, S \in \Omega^T} p(F, S | A_n) \\
&= \arg \max_{A_n \in A, S \in \Omega^T} [p(S | A_n) p(F | S, A_n)] \\
&= \arg \max_{A_n \in A, S \in \Omega^T} [p(S^1, \dots, S^T | A_n) p(F^1, \dots, F^T | S^1, \dots, S^T, A_n)] .
\end{aligned} \tag{2}$$

248 Assume 1) F is statistically independent of A_n given S , 2) F^t is statistically
249 dependent only on S^t , and 3) S^t is a Markov chain in an action, i.e., S^t only
250 depends on its previous state S^{t-1} . We can further simplify (2) as

$$\begin{aligned}
A^* &= \arg \max_{A_n \in A, S \in \Omega^T} \prod_{t=1}^T [p(S^t | S^{t-1}, A_n) p(F^t | S^t)] \\
&= \arg \max_{A_n \in A, S \in \Omega^T} \sum_{t=1}^T [\log(p(S^t | S^{t-1}, A_n)) + \log(p(F^t | S^t))] .
\end{aligned} \tag{3}$$

251 To solve (3), we adopt the Action-Specific Viterbi Decoding(ASVD) method [26]
252 and derive a new score function:

$$\begin{aligned}
&Score(A_n) \\
&= \max_{S \in \Omega^T} \sum_{t=1}^T [\log(p(S^t | S^{t-1}, A_n)) + \log(p(F^t | S^t))] \\
&= \max_{I \in \{1, \dots, K\}^T} \sum_{t=1}^T [\log(p(\omega_{I^t} | \omega_{I^{t-1}}, A_n)) + \log(p(F^t | \omega_{I^t}))] \\
&= \max_{I \in \{1, \dots, K\}^T} \sum_{t=1}^T [\log(B_n[I^{t-1}, I^t]) + \log(P_{I^t}(F^t))] \\
&= \max_{I \in \{1, \dots, K\}^T} \sum_{t=1}^T [\log(B_n[I^{t-1}, I^t]) - C \|F^t - \omega_{I^t}\|^2] ,
\end{aligned} \tag{4}$$

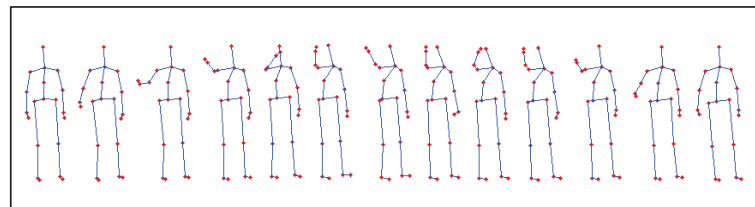
253 and

$$A^* = \arg \max_{A_n} Score(A_n) . \tag{5}$$

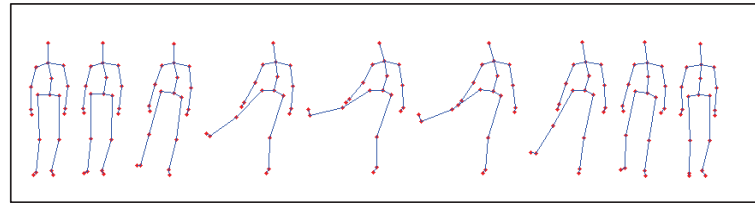
254 In (4), $I = [I^1, \dots, I^T]$ is a sequence of numbers, where $I^t \in \{1, \dots, K\}$,
 255 $\omega_{I^t} = S^t$, and $C = \frac{1}{2\sigma^2}$, which could be optimized by cross-validation using
 256 the training set if σ is hard to be reliably estimated.

257 4. Experiments and Discussion

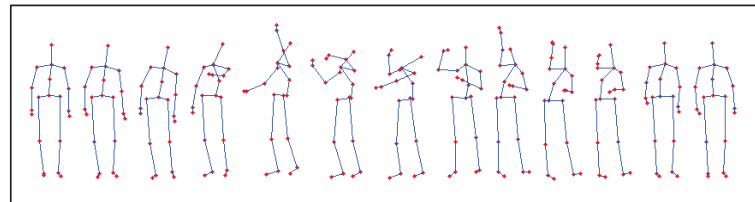
258 4.1. Datasets



(a) High Wave



(b) Side Kick



(c) Tennis Serve

Figure 7: Sample frames of 20-joints skeleton for actions of (a) High Wave, (b) Side Kick and (c) Tennis Serve from the MSRAction3D dataset.

259 **MSRAction3D dataset** [6]: The Microsoft Research Action3D dataset

260 (“MSRAction3D” for short) consists of 20 actions of 10 subjects, each action
261 with 2 or 3 repetitions. These actions are mainly interactions with console
262 in video games. As shown in Fig. 7, actions in this dataset capture a variety
263 of motions related to arms, legs, torso, and their combinations. Meanwhile,
264 the skeleton positions in this dataset are quite noisy. Hence, experiments
265 on this dataset were widely adopted to test the accuracy and robustness of
266 recognition methods for various actions. In previous works, it was used in
267 two ways, either as one dataset containing all actions [7, 9–11, 14, 19–21] or
268 by division into three subsets [6, 8, 12, 13, 15–19, 21].

269 **UTKinect-Action dataset** [13]: The University of Texas Kinect Action
270 dataset (“UTKinect-Action” for short) consists of 10 actions of 10 different
271 persons, each action with 1 or 2 repetitions. These actions, including *walk*,
272 *sit down*, *stand up*, *pick up*, *carry*, *throw*, *push*, *pull*, *wave* and *clap hands*,
273 are mainly obtained from daily life. One of the challenges of this dataset is
274 that one of the persons is left-handed. Meanwhile, the lengths of samples
275 from this dataset vary in a wide range. It was used for evaluation in several
276 previous works [13, 16, 17, 19].

277 4.2. Experiments and Results

278 In all experiments, we use the same hardware setup: Intel Core i7-4790
279 CPU @3.6GHz, 24GB RAM. Our method could run at about 30fps for test.

280 In all experiments, C in (4) and η in (1) are tuned by cross-validation
281 using the training set, which follows the procedure below: (i) set a grid of
282 values of C and η , divide the training set into two halves; (ii) train on the
283 first half with different combinations of C and η , get accuracies by testing
284 on the second half; (iii) repeat (ii) but exchange the two halves, average the

285 accuracies from (ii) and (iii); and (iv) find the values of C and η that produce
286 the highest accuracy.

287 4.2.1. *MSRAAction3D*

288 We follow the two types of experiments, as with the previous works.
289 Firstly, the whole dataset is used to verify the performance of our method on
290 a large number of actions. Secondly, three subsets of this dataset are tested
291 to confirm the applicability of the method in various situations.

292 In the first type of experiments, the entire 20 actions with 557 sequences
293 are applied as with [9]. We conduct the cross-subject test, where sequences of
294 half the subjects are used for training and the rest for testing. We repeat the
295 experiment 252 times for different folds of 10 subjects as with [9]. The per-
296 formance, including the best result, the worst result and the average result,
297 are listed and compared in Table 1. The confusion matrix of the best result
298 is displayed in Fig. 8. In the confusion matrix, the vertical coordinate (y)
299 represents the true label of an action sequence and the horizontal coordinate
300 (x) represents the recognition result. The value at the (x, y) coordinate of
301 the matrix represents the ratio of action y recognized as action x . As shown
302 in Table 1, our method outperforms other methods in terms of the average
303 result and/or the best result. From the confusion matrix of the best result
304 in Fig. 8, we can see that 15 out of 20 actions achieve 100% accuracy, which
305 is perfect considering the noise of skeleton collected by Kinect and the huge
306 intra-class variance among different subjects. Note that the action *Hand*
307 *Catch* gets the lowest accuracy 50%, mainly because it is similar to other
308 two actions *High Wave* and *High Throw*.

309 In the second type of experiments, as with [6] all the 567 sequences of the

Table 1: Recognition rates of the first experiment on MSRAction3D. In the table, we present the best result, the worst result and the average result for 252 different folds of 10 subjects.

Method	Avg \pm Std	Best	Worst
Ours	87.05 \pm 3.75	95.56	74.39
Random Occupancy Patterns [7]	-	86.50	-
Actionlet Ensemble [11]	-	88.2	-
HON4D+ D_{disc} [9]	82.15 \pm 4.18	88.89	-
Spatial and Temporal Part Sets [14]	-	90.22	-
HOPC of 3D Pointclouds [10]	86.49 \pm 2.28	92.39	74.36
Points in a Lie Group [19]	-	89.48	-
Histograms of Action Poses + DTW [20]	-	90.56	-
Dynemes and Forward Differences [21]	-	91.94	-

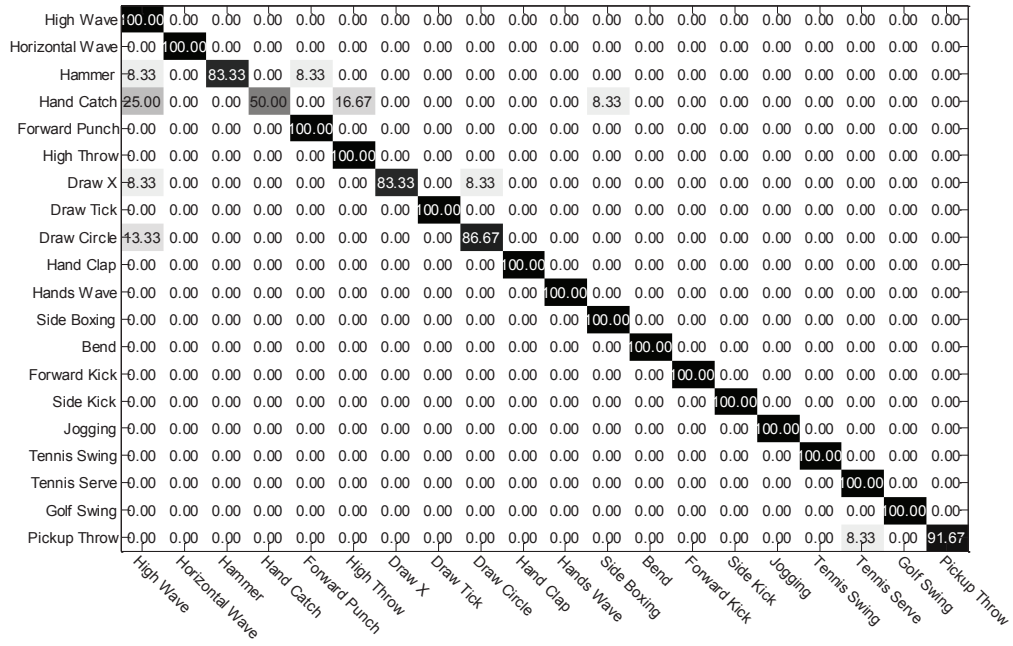


Figure 8: Confusion matrix of our method on MSRAction3D with the part-based clustering. In the confusion matrix, the vertical coordinate (y) represents the true label of an action sequence and the horizontal coordinate (x) represents the recognition result. The value at the (x, y) coordinate of the matrix represents the ratio of action y recognized as action x.

310 dataset are split into three subsets, each with eight actions. Evaluation for
 311 each subset is done independently. Table 2 shows the three subsets in this
 312 experiment. These three subsets have different aims: AS1 and AS2 contain
 313 similar actions to verify a method’s ability to discriminate similar movements,
 314 while AS3 groups complex actions together to evaluate the versatility of a
 315 method. The overall recognition rate is calculated by averaging the results
 316 over subsets. We also use the cross-subject test for this experiment and
 317 repeat it for different folds of subjects. The results are shown in Table 3,
 318 from which we can observe that our method outperforms the state-of-the-art
 319 method [21] by 2.5% in terms of the best result.

Table 2: The three subsets of actions for the second experiment on MSRAction3D.

Action Set 1 (AS1)	Action Set 2 (AS2)	Action Set 3 (AS3)
Horizontal Wave (HoW)	High Wave (HiW)	High Throw (HT)
Hammer (H)	Hand Catch (HC)	Forward Kick (FK)
Forward Punch (FP)	Draw X (DX)	Side Kick (SK)
High Throw (HT)	Draw Tick (DT)	Jogging (J)
Hand Clap (HC)	Draw Circle (DC)	Tennis Swing (TSw)
Bend (B)	Hands Wave (HW)	Tennis Serve (TSr)
Tennis Serve (TSr)	Forward Kick (FK)	Golf Swing (GS)
Pickup Throw (PT)	Side Boxing (SB)	Pickup Throw (PT)

Table 3: Recognition rates of the second experiment on MSRAction3D.

Method	Avg \pm Std	Best	Worst
Ours	88.7 \pm 3.6	96.1	78.0
Bag of 3D Points [6]	-	74.7	-
Histograms of 3D Joints [13]	-	78.97	-
EigenJoints [12]	-	82.3	-
EigenJoints + Hierarchical [18]	-	90.3	-
Histograms of Oriented Displacements [15]	-	91.26	-
Random Forests [16]	-	94.3	-
Space-Time Pose [17]	-	92.77	-
Points in a Lie Group [19]	-	92.46	-
Dynemes and Forward Differences [21]	-	93.6	-

320 *4.2.2. UTKinect-Action*

321 We follow the challenging cross-subject test setting of [16, 19] instead of
 322 leave-one-out-cross-validation (LOOCV) setting of [13, 17]. In this setting,
 323 half of the subjects are used for training while the remaining for testing. The
 324 performances on this dataset are compared in Table 4. We can find out that
 325 our method, comparable to the method of Vemulapalli *et al.* [19], is much
 326 better than other methods. The confusion matrix of our method is shown in
 327 Fig.9, from which we can observe that all 10 actions achieve accuracy higher
 328 than 90%, 6 out of which achieve 100% accuracy.

Table 4: Recognition rates on UTKinect-Action.

Method	Recognition rate
Ours	95.96
Histograms of 3D Joints [13]	90.92
Random Forest [16]	91.9
Space-Time Pose [17]	91.5
Points in a Lie Group [19]	97.08

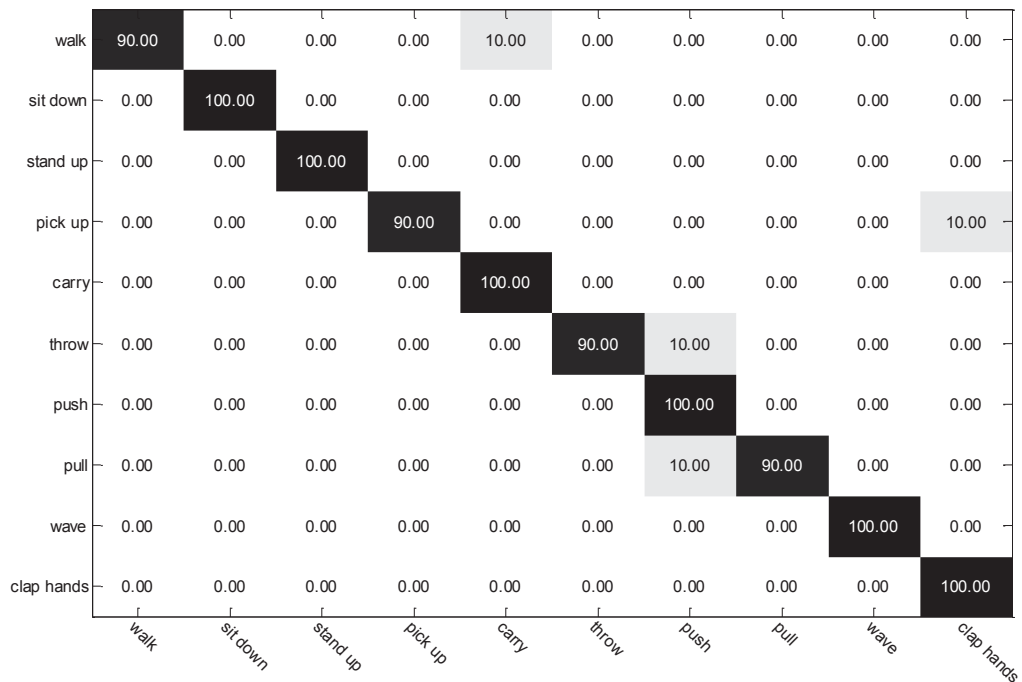


Figure 9: Confusion matrix of our method on UTKinect-Action.

329 4.3. Module Analysis

330 We carry out further experiments to analyze the effects of the three mod-
331 ules. All these experiments follow the same setup as the first type of ex-
332 periments on MSRAction3D aforementioned, i.e. the entire 20 actions being
333 used, cross-subject test and repeating for different folds of subjects.

334 4.3.1. Part-based Clustering

335 We perform the first type of experiments on MSRAction3D to verify the
336 part-based clustering. The results are shown in Table 5. In terms of the av-
337 erage result, the part-based clustering significantly improves the recognition
338 rate by about 14%. The improvement is mainly due to two merits of the
339 part-based clustering: 1) it automatically divides distinct sequences of the
340 same action into more than one cluster, and 2) it provides the relevant joints
341 for motion feature extraction to help distinguish similar actions in the same
342 cluster. Fig. 5 shows the part-based clustering result obtained from the 20
343 actions.

344 For the first merit, again take action *Side Boxing* as an example. Some-
345 body may do it with one hand while others may prefer to use two hands,
346 which results in a large intra-class variance (see Fig. 1). Without the part-
347 based clustering, it only achieves 86.7% recognition accuracy. If we restrict
348 that the sequences of the same action could only be grouped into a single
349 cluster, it is difficult to decide which cluster, the cluster of using one hand
350 or the cluster of using two hands, is correct for *Side Boxing*. Using either
351 one may lead to misclassifying the other kind of sequences, which will result
352 in a low recognition accuracy of 70%, even lower than the result without
353 the part-based clustering. In contrast, allowing them to appear in different

354 clusters (see Fig. 5) can solve the problem, which leads to 100% recognition
 355 accuracy for this action in our experiment. A similar case is with action
 356 *Jogging*.

357 For the second merit, take actions *Hammer*, *High Throw* and *Draw Cir-*
 358 *cle* as an example. Without the part-based clustering, action *Hammer* is
 359 highly confused with *High Throw* or *Draw Circle* and only 25% are recog-
 360 nized correctly. As these actions are similar movements using only the right
 361 hand, motion features extracted from all joints of body would contain much
 362 irrelevant information as noise, resulting in a low recognition accuracy. With
 363 the part-based clustering, they are grouped into the same cluster and only
 364 the joints of right up extreme (RUE) are used for motion feature extrac-
 365 tion. These motion features are the most discriminative with the irrelevant
 366 noise dropped out, which ultimately improves the recognition rate of *Ham-*
 367 *mer* up to 83.3%. Similar situations exist with *Horizontal Wave* (confused
 368 with *Draw X* and *Draw Tick*), *Draw X* (confused with *Draw Circle*) and
 369 *Jogging* (confused with *Forward Kick*).

Table 5: Performance comparison between methods with and without the part-based clustering.

Method	Avg \pm Std	Best	Worst
With part-based clustering	87.05 \pm 3.75	95.56	74.39
Without part-based clustering	73.08 \pm 4.34	83.03	61.97

Table 6: Performance comparison among different time scales for feature extraction. The notation 1, 5, 10 means that three previous frames (1, 5, 10) before the current frame are used.

Time scale	Avg \pm Std	Best	Worst
1	81.51 \pm 3.66	89.59	70.38
1, 5	84.88 \pm 3.35	92.67	75.61
1, 5, 10	87.05 \pm 3.75	95.56	74.39
1, 5, 10, 15	82.31 \pm 4.00	93.70	72.47
1, 5, 10, 15, 20	79.78 \pm 4.28	93.70	68.77

370 *4.3.2. Time Scale for Motion Feature Extraction*

371 As described in Section 3.2, to take the time scale into consideration, we
372 calculate the *DP* component of the motion features with several previous
373 frames before the current frame. To determine how many of previous frames
374 are sufficient, we repeat the first type of experiments on MSRAction3D with
375 different numbers of previous frames for motion feature extraction and com-
376 pare the recognition accuracy in Table 6. We can observe that taking the
377 time scale into consideration could enhance representability of the motion
378 features, which in turn results in higher recognition accuracy. Besides, us-
379 ing three previous frames (1, 5, 10) provides the best performance in this
380 experiment. Take the action *Bend* as an example. Without taking the time
381 scale into consideration, it only achieves 58.3% recognition rate. While using
382 two previous frames (1, 5) already improves it up to 91.7% and using three
383 previous frames (1, 5, 10) even manages to achieve 100% recognition rate.
384 However, using more previous frames (15, 20) leads to worse performance,

385 which is mainly because the previous frames far away could not provide valu-
 386 able information for recognition but produce more noise.

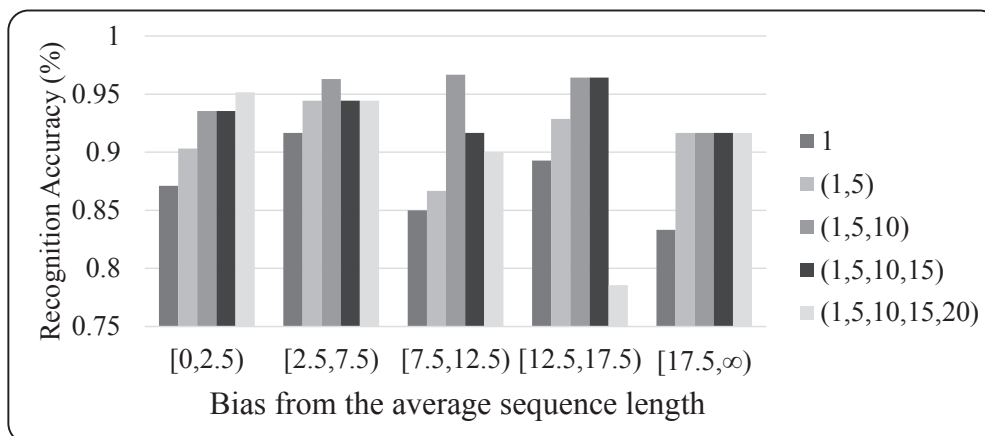


Figure 10: Recognition accuracies of various sequence lengths using different time scales for motion feature extraction. In the figure, “Bias from the average sequence length” means the absolute difference between the length of a sequence and the average length of all sequences of this action, according to which we divide the test set into several groups and compare the recognition accuracies for each group using different time scales.

387 In order to investigate how to decide the time scale for a given dataset,
 388 we record the sequence lengths of the samples of all 20 actions. We find out
 389 that for this dataset, the standard deviations of sequence lengths of most
 390 actions are around 10, which explains why using up to 10 previous frames is
 391 sufficient to represent the time scale. We further partition the test set into
 392 several groups of sequences according to the difference between the length of
 393 a sequence and the average sequence length, that is, according to the bias
 394 from the average sequence length. For each group, we then compare the

395 recognition accuracies obtained from using different time scales for motion
396 feature extraction. The comparative results are plotted in Fig. 10 (accuracies
397 are obtained from the best result). We can observe that extracting motion
398 features for three previous frames (1, 5, 10) improves the accuracies for differ-
399 ent sequence lengths. Moreover, the improvement is relatively prominent for
400 sequences whose length biases from the average are around 10 ([7.5, 12.5]).
401 These observations imply that, given a dataset, using the time scale up to the
402 standard deviations of sequence lengths is highly possible to offer an optimal
403 performance.

404 *4.3.3. Action Graphs*

405 The Viterbi decoding algorithm used for action graphs is a dynamic pro-
406 gramming algorithm, which could output scores (confidences) for all actions
407 at any frame. This suggests that we may make use of this benefit to recognize
408 actions at earlier time before the end of an action with sufficient confidence or
409 to automatically and reliably segment actions along with action recognition.
410 This is also useful to make the recognition latency lower when we apply the
411 recognition algorithm in real-time applications.

412 Here we carry out an experiment to verify the early detection performance
413 of our method. In this experiment, we rescale the length of all test action
414 sequences to be 1, and recognize the action at the points of $5/6$, $2/3$, $1/2$,
415 $5/12$, $1/3$, $1/4$, respectively. (Here we do not set the points to make equal
416 segments, just because the recognition rate usually does not vary equably.
417 Thus we sample more points for the early stage where the recognition rate
418 varies more rapidly.) The training procedure is just the same as before. We
419 record the recognition rates at different points for all folds of subjects and

420 put the average results and the standard deviation together to draw Fig. 11.
 421 We can observe that, from right to left, at first the recognition rate goes down
 422 slowly with the recognition point moving earlier, and then it drops quickly
 423 when the point is earlier than 2/3. So with bearable decline of recognition
 424 accuracy (about 4%), we could achieve early detection of actions at the 2/3
 425 point of the action sequence. At this point, the performance of our method
 426 is still higher than [9] (see Table 1).

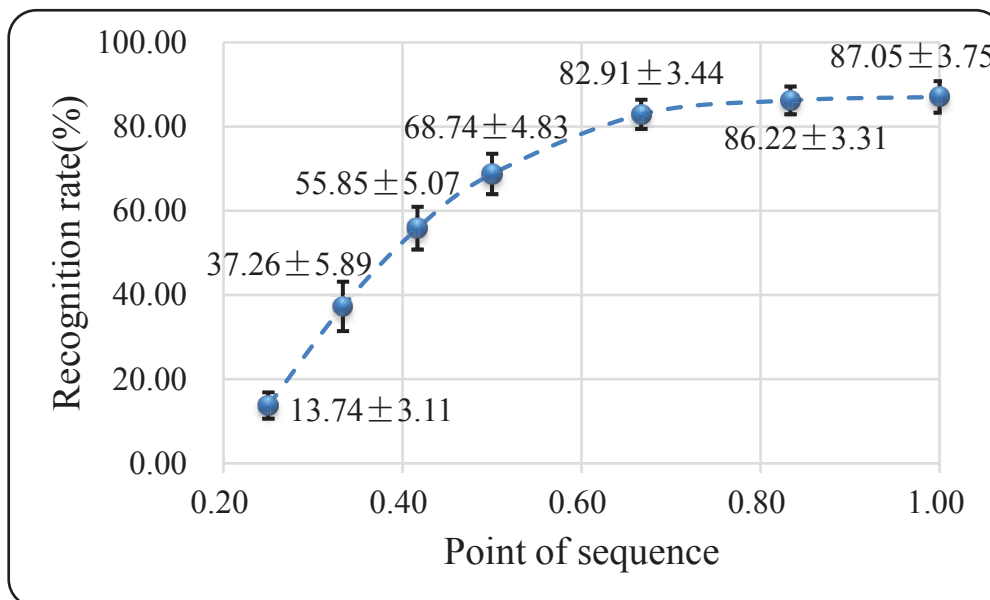


Figure 11: Recognition rates at different points of action sequences. We rescale the length of all test action sequences to be 1, and a point means the length used for prediction.

427 To summarize up the module analysis, the introduction of the part-based
 428 clustering in our proposal and modifying the motion features of [12] with con-
 429 sideration of time scale are demonstrated to be effective on the MSRAction3D

430 dataset, resulting in better performance. Meanwhile, we suggest a statistical
431 principle for deciding the time scale for any specified dataset. Furthermore,
432 early detection is proved to be feasible with the use of action graphs.

433 *4.4. Impact of Parameters*

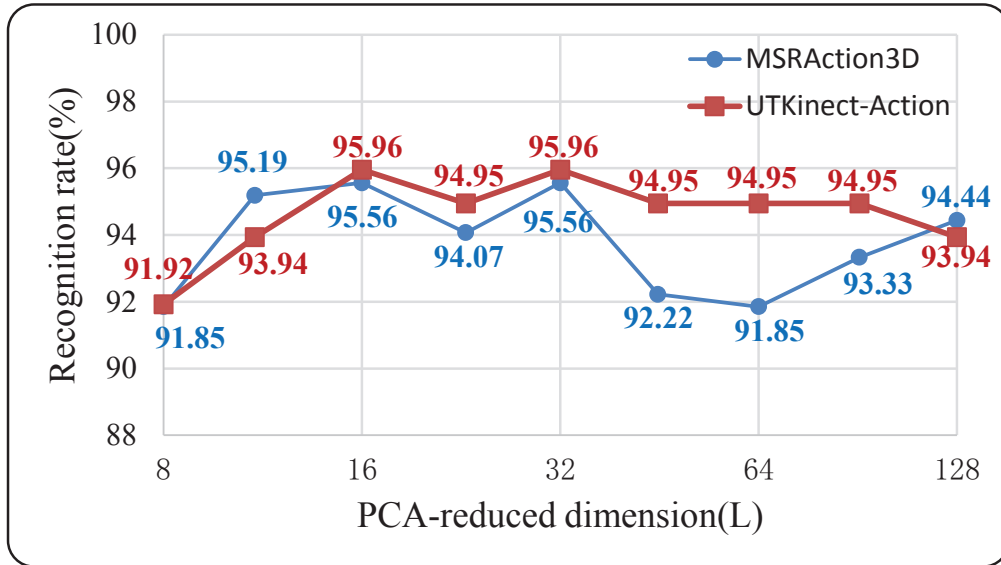
434 To verify the performance of our method versus different parameter values
435 (here parameters include the PCA-reduced dimension L and the posture
436 number K), we carry out experiments on MSRAction3D (the first type) and
437 UTKinect-Action with various parameter values for the fold that gets the
438 best result.

439 *4.4.1. Impact of the PCA-reduced Dimension L*

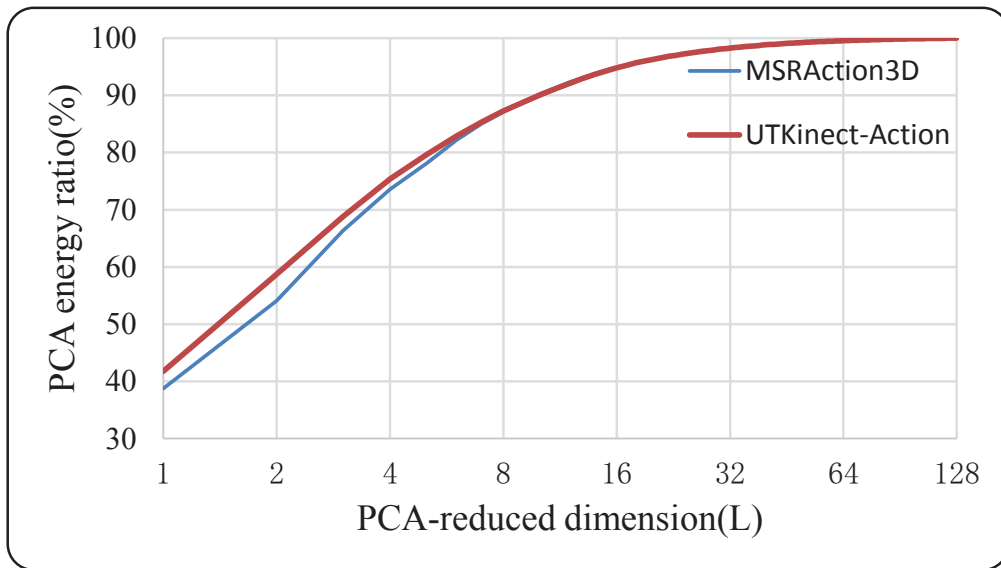
440 Fig. 12a shows how the performance of our method varies with different
441 values of the PCA-reduced dimension, L . We can observe that, for both
442 datasets, when L is small the performance slightly increases with L , and
443 when L is large enough the performance becomes rather stable and then
444 slightly decreases. This may be because when L is sufficiently large we can
445 capture most of the energy as shown in the energy graph, Fig. 12b, and then
446 more unimportant principal components induce certain irrelevant informa-
447 tion. When L reaches 16, we already capture 95% of the energy and achieve
448 the best performance. For this reason and for the consideration of small
449 computational costs, we choose 16 as the value of L for both datasets.

450 *4.4.2. Impact of the Posture Number K*

451 Fig. 13 shows the impact of the number of postures, K , on the recognition
452 performance. As we need more postures to characterize more actions, we take
453 ratio $\rho = K/N$ as the parameter to be discussed, where N is the number of



(a) Impact of the PCA-reduced dimension on recognition performance.



(b) PCA energy graph. In the graph, “PCA energy ratio” is defined as $\frac{\sum_{i=1}^L \lambda_i}{\sum_i \lambda_i}$, where λ_i are the eigenvalues obtained from PCA and they are ordered as $\lambda_i \geq \lambda_{i+1}$.

Figure 12: Impact of the PCA-reduced dimension.

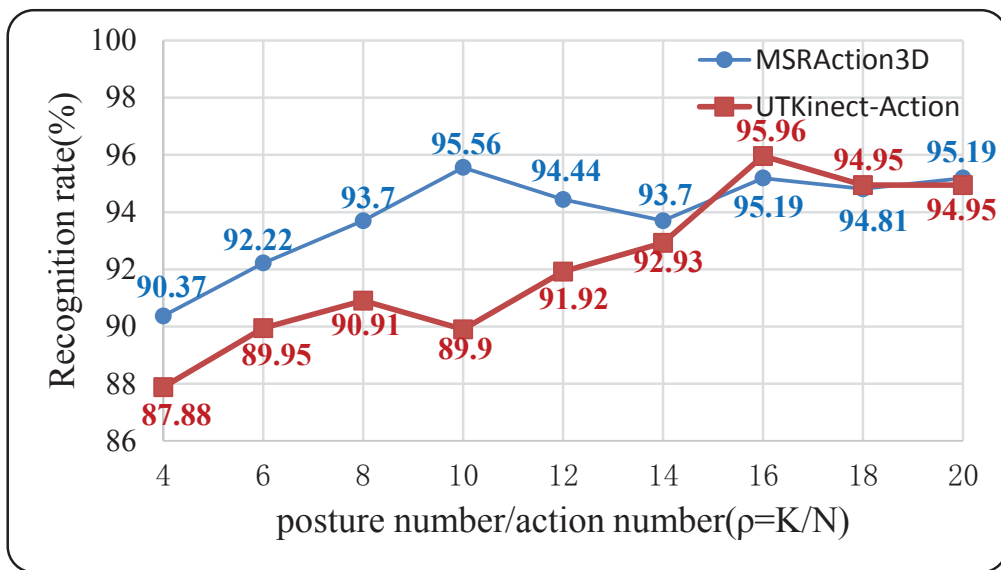


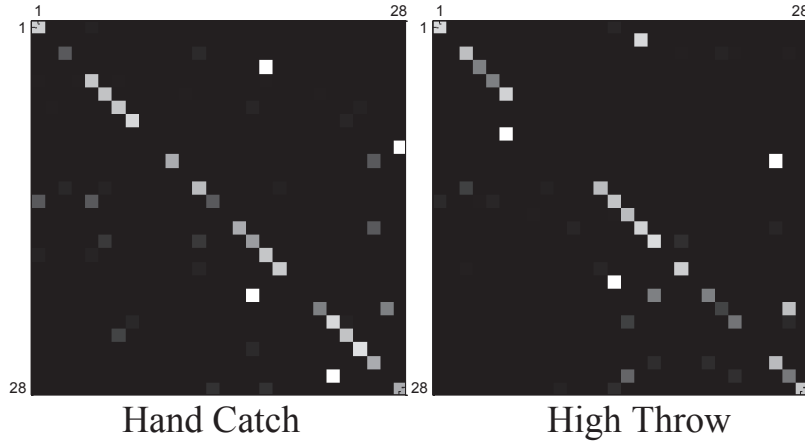
Figure 13: Impact of ρ on the performance. Here $\rho = K/N$, where N is the number of actions and K is the number of postures.

454 actions. As we can observe in Fig. 13, as ρ increases, the performance at first
 455 improves significantly and then stays at a high level.

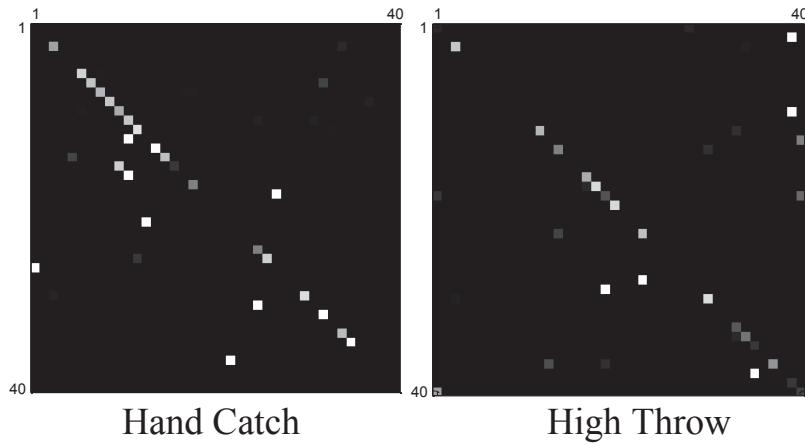
456 To explain this pattern, take two actions *Hand Catch* and *High Throw*
 457 in the MSRAction3D dataset as an example. When ρ is too small, we do
 458 not have enough postures to characterize all the actions so the transitional
 459 probability matrices of them are similar, which results in a poor recognition
 460 performance. Fig. 14a shows the main part of the transitional probab-
 461 ity matrices of the two actions when ρ is too small (4). The two matrices
 462 seem similar to each other, which leads to a low recognition rate of 25% for
 463 *Hand Catch*. As ρ gets large enough, we manage to characterize different
 464 actions with different postures and/or different transitional probability ma-
 465 trices, which makes the recognition task much easier. Fig. 14b shows the
 466 main part of the transitional probability matrices of the same two actions
 467 when ρ is large enough (10). Compared with those in Fig. 14a where ρ is too
 468 small, the two matrices now are much more different from each other, which
 469 improves the recognition rate of *Hand Catch* to 75%. To verify the visual
 470 impression of the matrix similarity in Fig. 14a and Fig. 14b, we calculate the
 471 correlation r of the two transitional probability matrices, A and B , for these
 472 two actions, using the following formula:

$$r = \frac{\sum_{i=1}^N \sum_{j=1}^N (A_{i,j} - \bar{A})(B_{i,j} - \bar{B})}{\sqrt{\left(\sum_{i=1}^N \sum_{j=1}^N (A_{i,j} - \bar{A})^2\right) \left(\sum_{i=1}^N \sum_{j=1}^N (B_{i,j} - \bar{B})^2\right)}}$$

473 where $\bar{A} = \frac{\sum_{i=1}^N \sum_{j=1}^N A_{i,j}}{N^2}$, $\bar{B} = \frac{\sum_{i=1}^N \sum_{j=1}^N B_{i,j}}{N^2}$. When ρ is 4, r is 0.44; when
 474 ρ is 10, r is 0.20, which indicates that the matrices are much less similar.
 475 Considering Fig. 13 and reasonable computational costs, we set ρ to 10 for



(a) When ρ is too small (4), matrices of two actions are moderately similar ($r = 0.44$).



(b) When ρ is large enough (10), two matrices are not similar ($r = 0.20$).

Figure 14: Transitional probability matrices of two actions at different values of ρ . The value at the (i, j) coordinate of the matrices represents the transitional probabilities from posture i to posture j . Each value is illustrated by the brightness, with a brighter one for a larger probability.

476 the MSRAction3D dataset and 16 for the UTKinect-Action dataset.

477 In summary, all the above experiments have verified that our proposed
478 method has tackled to some extent the three challenges mentioned at the
479 beginning of this paper. Firstly, with the utilization of the part-based clus-
480 tering module, the challenge of high intra-class variance with low inter-class
481 variance has been mostly solved. Secondly, we have notably worked out the
482 challenge of variable movement speed by considering time scales for motion
483 features. Thirdly, since our method is based on 3D skeleton, its computation
484 costs are relatively low, making the method applicable in real time.

485 5. Conclusion

486 In this paper, we have proposed a novel two-level hierarchical framework
487 for action recognition with 3D skeleton sequences. In the framework, we
488 have introduced a new part-based five-dimensional feature vector to mine
489 the most relevant body parts for each action sequence and to cluster action
490 sequences, have investigated the time scale of dynamics to optimally modify
491 established motion features, and have devised a score function for the action
492 inference based on action graphs. Our experiments have also verified that,
493 compared with other state-of-the-art methods, the proposed method could
494 achieve higher accuracy on the complex MSRAction3D dataset.

495 Nevertheless, the performance of our method still considerably depends
496 on the accuracy of the skeleton positions, even though we only utilize the
497 relevant joints. Besides, the two datasets used here provide known beginning
498 and end of the actions, which are not available in real-world interactions.
499 Hence, further to our work is to investigate the open problem of segment-

500 ing actions automatically, reliably and quickly by using action graphs, as
501 suggested in Section 4.3.3.

502 **Acknowledgements**

503 The authors are grateful to two reviewers and Mr Hengkai Guo for their
504 constructive comments, in particular for their suggestions which have led to
505 a large improvement of section 4.3, section 4.4 and the readability of this
506 paper. Thanks to Mr Andrew Mpapalika for proofreading. This work was
507 partially supported by NSFC61271390 and 2015AA016304.

508 **References**

- 509 [1] R. Poppe, A survey on vision-based human action recognition, *Image*
510 *and Vision Computing* 28 (6) (2010) 976–990.
- 511 [2] G. Wang, X. Yin, X. Pei, C. Shi, Depth estimation for speckle projec-
512 tion system using progressive reliable points growing matching, *Applied*
513 *Optics* 52 (3) (2013) 516–524.
- 514 [3] X. Yin, G. Wang, C. Zhang, Q. Liao, Learning the missing values in
515 depth maps, in: *International Conference on Optical Instruments and*
516 *Technology (OIT2013)*, International Society for Optics and Photonics,
517 2013, pp. 904508–904508.
- 518 [4] K. Liu, C. Zhou, S. Wei, S. Wang, X. Fan, J. Ma, Optimized stereo
519 matching in binocular three-dimensional measurement system using
520 structured light, *Applied Optics* 53 (26) (2014) 6083–6090.

- 521 [5] C. Shi, G. Wang, X. Yin, X. Pei, B. He, X. Lin, High-accuracy stereo
522 matching based on adaptive ground control points, *Image Processing,*
523 *IEEE Transactions on* 24 (4) (2015) 1412–1423.
- 524 [6] W. Li, Z. Zhang, Z. Liu, Action recognition based on a bag of 3D points,
525 in: *Computer Vision and Pattern Recognition Workshops (CVPRW),*
526 *2010 IEEE Computer Society Conference on,* IEEE, 2010, pp. 9–14.
- 527 [7] J. Wang, Z. Liu, J. Chorowski, Z. Chen, Y. Wu, Robust 3D action
528 recognition with random occupancy patterns, in: *Computer Vision–*
529 *ECCV 2012,* Springer, 2012, pp. 872–885.
- 530 [8] X. Yang, C. Zhang, Y. Tian, Recognizing actions using depth motion
531 maps-based histograms of oriented gradients, in: *Proceedings of the*
532 *20th ACM International Conference on Multimedia,* ACM, 2012, pp.
533 1057–1060.
- 534 [9] O. Oreifej, Z. Liu, HON4D: Histogram of oriented 4D normals for activ-
535 ity recognition from depth sequences, in: *Computer Vision and Pattern*
536 *Recognition (CVPR), 2013 IEEE Conference on,* IEEE, 2013, pp. 716–
537 723.
- 538 [10] H. Rahmani, A. Mahmood, D. Q Huynh, A. Mian, HOPC: Histogram of
539 oriented principal components of 3D pointclouds for action recognition,
540 in: *Computer Vision–ECCV 2014,* Springer, 2014, pp. 742–757.
- 541 [11] J. Wang, Z. Liu, Y. Wu, J. Yuan, Mining actionlet ensemble for ac-
542 tion recognition with depth cameras, in: *Computer Vision and Pattern*

- 543 Recognition (CVPR), 2012 IEEE Conference on, IEEE, 2012, pp. 1290–
544 1297.
- 545 [12] X. Yang, Y. Tian, Eigenjoints-based action recognition using naive-
546 Bayes-nearest-neighbor, in: Computer Vision and Pattern Recognition
547 Workshops (CVPRW), 2012 IEEE Computer Society Conference on,
548 IEEE, 2012, pp. 14–19.
- 549 [13] L. Xia, C.-C. Chen, J. Aggarwal, View invariant human action recog-
550 nition using histograms of 3D joints, in: Computer Vision and Pattern
551 Recognition Workshops (CVPRW), 2012 IEEE Computer Society Con-
552 ference on, IEEE, 2012, pp. 20–27.
- 553 [14] C. Wang, Y. Wang, A. L. Yuille, An approach to pose-based action
554 recognition, in: Computer Vision and Pattern Recognition (CVPR),
555 2013 IEEE Conference on, IEEE, 2013, pp. 915–922.
- 556 [15] M. A. Gowayyed, M. Torki, M. E. Hussein, M. El-Saban, Histogram of
557 oriented displacements (HOD): describing trajectories of human joints
558 for action recognition, in: Proceedings of the Twenty-Third interna-
559 tional joint conference on Artificial Intelligence, AAAI Press, 2013, pp.
560 1351–1357.
- 561 [16] Y. Zhu, W. Chen, G. Guo, Fusing spatiotemporal features and joints for
562 3D action recognition, in: Computer Vision and Pattern Recognition
563 Workshops (CVPRW), 2013 IEEE Conference on, IEEE, 2013, pp. 486–
564 491.

- 565 [17] M. Devanne, H. Wannous, S. Berretti, P. Pala, M. Daoudi,
566 A. Del Bimbo, Space-time pose representation for 3D human action
567 recognition, in: *New Trends in Image Analysis and Processing–ICIAP*
568 2013, Springer, 2013, pp. 456–464.
- 569 [18] H. Chen, G. Wang, L. He, Accurate and real-time human action recog-
570 nition based on 3D skeleton, in: *International Conference on Optical*
571 *Instruments and Technology (OIT2013)*, International Society for Op-
572 tics and Photonics, 2013, pp. 90451Q–90451Q.
- 573 [19] R. Vemulapalli, F. Arrate, R. Chellappa, Human action recognition by
574 representing 3D skeletons as points in a Lie group, in: *Computer Vision*
575 *and Pattern Recognition (CVPR)*, 2014 IEEE Conference on, IEEE,
576 2014, pp. 588–595.
- 577 [20] M. Barnachon, S. Bouakaz, B. Boufama, E. Guillou, Ongoing human ac-
578 tion recognition with motion capture, *Pattern Recognition* 47 (1) (2014)
579 238–247.
- 580 [21] I. Kapsouras, N. Nikolaidis, Action recognition on motion capture data
581 using a dynemes and forward differences representation, *Journal of Vi-*
582 *sual Communication and Image Representation* 25 (6) (2014) 1432–1445.
- 583 [22] F. Offi, R. Chaudhry, G. Kurillo, R. Vidal, R. Bajcsy, Sequence of the
584 most informative joints (SMIJ): A new representation for human skeletal
585 action recognition, *Journal of Visual Communication and Image Repre-*
586 *sentation* 25 (1) (2014) 24–38.

- 587 [23] J. Shotton, R. Girshick, A. Fitzgibbon, T. Sharp, M. Cook, M. Finoc-
588 chio, R. Moore, P. Kohli, A. Criminisi, A. Kipman, A. Blake, Efficient
589 human pose estimation from single depth images, *Pattern Analysis and*
590 *Machine Intelligence, IEEE Transactions on* 35 (12) (2013) 2821–2840.
- 591 [24] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio,
592 A. Blake, M. Cook, R. Moore, Real-time human pose recognition in
593 parts from single depth images, *Communications of the ACM* 56 (1)
594 (2013) 116–124.
- 595 [25] L. He, G. Wang, Q. Liao, J.-H. Xue, Depth-images-based pose estima-
596 tion using regression forests and graphical models, *Neurocomputing* 164
597 (2015) 210–219.
- 598 [26] W. Li, Z. Zhang, Z. Liu, Expandable data-driven graphical modeling
599 of human actions based on salient postures, *Circuits and Systems for*
600 *Video Technology, IEEE Transactions on* 18 (11) (2008) 1499–1510.