

Adaptive Segmentation and Sequence Learning of Human Activities from Skeleton Data

David Ada Adama, Ahmad Lotfi*, Robert Ranson

School of Science and Technology, Nottingham Trent University, Clifton Lane, Nottingham, NG11 8NS, United Kingdom

Abstract

Discovering underlying patterns for predicting future actions from spatio-temporal human activity information is a fundamental component of research related to the development of expert systems in human activity recognition and assistive robotics. Current research focuses on classification or learning representations of activities for various applications. However, not much attention is given to the pattern discovery of activities which have a major role in the prediction of unseen actions. This paper proposes a novel Adaptive Segmentation and Sequence Learning (ASSL) framework which aims at segmenting unlabelled observations of human activities from extracted 3D joint information. Learning from these obtained segments provides information about the underlying patterns of activity sequences needed in predicting subsequent actions. In the proposed method, the temporal accumulated motion energy of body parts in an activity is utilised in the segmentation process to obtain key actions from unlabelled activity sequences since body parts show changes in acceleration and deceleration during an activity. Based on the segments obtained, the temporal sequence of transitions across activity segments are learned by employing a Long Short-Term Memory Recurrent Neural Network. This ASSL technique has been evaluated using both an experimental human activity dataset and a public activity dataset, and achieved a better performance when compared with other techniques including an Auto-regressive Integrated Moving Average, Support Vector Regression and Gaussian Mixture Regression Models in learning to predict patterns of activity sequences.

*Corresponding author

Email addresses: david.adama2015@my.ntu.ac.uk (David Ada Adama), ahmad.lotfi@ntu.ac.uk (Ahmad Lotfi), robert.ranson@ntu.ac.uk (Robert Ranson)

Keywords: Human Activity Recognition, Sequence Learning, Long Short-Term Memory, Activity Segmentation, Mean-shift clustering, Key Pose.

1 **1. Introduction**

2 The advances in technology have seen more research on the development of expert
3 systems related to human activities and their applications in everyday life. Learning the
4 sequences of human activities is one aspect that is daunting in many of such applications.
5 Due to the variability in the human nature of conducting activities, it is often not possible
6 to attain a generalised model for identifying sequences used for activity predictions. This
7 is due to understanding the underlying patterns of activities which in many cases are not
8 explicit.

9 A popular area of the application of expert systems for human activity sequence learning
10 is human-robot interaction. For example, assistive robots require abilities to learn human
11 activities in order to function autonomously. Such activities usually require the coordination
12 of different joints in the body to accomplish activities such as “pick and place” of an object
13 activities. Robots equipped with preset instructions (or models) to carry out predefined
14 functions limits them to only certain tasks as they do not possess the intelligence required to
15 evolve their knowledge into executing functions which may differ from the preset knowledge.
16 Also, such models become obsolete as new tasks are encountered since they are not able
17 to adapt to dynamic situations which are inherent in most practical applications. This is
18 primarily due to variations in activity sequences, thus the need to investigate the varying
19 patterns of human activities. To offer a solution to such cases, it is imperative to break
20 down these activities into constituent elements and extract relevant information used in
21 simplifying the process of recognising various human activity patterns. Fig. 1 illustrates
22 the underlying concept of how human activity patterns can be inferred and learned from
23 processing extracted visual 3D information.

24 There are two main categories of learning algorithms suitable for human activity learning:
25 *Batch learning* and *Sequence learning*. Classical batch learning algorithms predict output for
26 new data when a complete training set of data is used. In this case, the new data samples are
27 presented simultaneously and when desired. However, a complete training dataset is often
28 not available in advance for most practical applications. In applications such as human
29 activity prediction (Adama et al., 2018), healthcare monitoring (Anderez et al., 2020) and

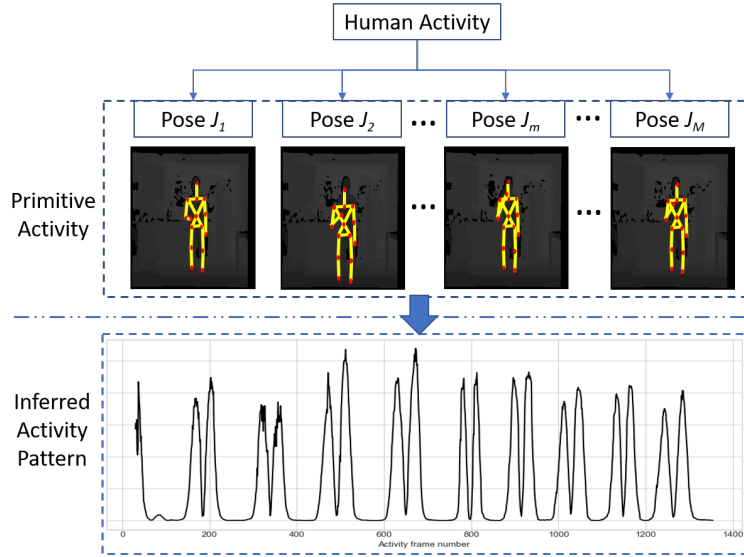


Fig. 1. An illustration of learning underlying patterns of simple primitive human activity sequences from 3D temporal information.

30 industrial functions (Suresh et al., 2010) in which temporal changes within a task are being
 31 observed, the classical batch learning algorithms are rather infeasible for learning. Sequence
 32 learning is executed in a series of occurrences of samples within a given training dataset.
 33 Samples are used in the algorithm one after another and discarded after learning. This
 34 implies that the computational time and memory required for learning is reduced, and the
 35 learning process can accommodate temporal changes associated with tasks (Suresh et al.,
 36 2010). In most cases of humans executing tasks, the path of actions may vary, however,
 37 each path contains approximately a similar order of true segments. To effectively learn
 38 such sequences of tasks, there are two key challenges which are often encountered. The
 39 *segmentation* of tasks wherein given the observed task path, the start and end positions
 40 of constituent actions through the path are identified and *sequential learning* of essential
 41 underlying actions (Lioutikov et al., 2017). The task segmentation is critical in sequence
 42 learning for modelling and interpreting tasks information as it facilitates the adaptation of
 43 learning sequences in unseen situations (Krishnan et al., 2017).

44 The main contributions of the work presented in this paper are summarised as follows:

- 45 - The paper proposes a novel adaptive segmentation and sequence learning (ASSL)
 46 approach for human activity pattern discovery from unlabelled sequences of observed

47 activities.

- 48 - Exploiting the temporal accumulated motion energy of human actions through activity
49 sequences for extracting key actions points during activities.
- 50 - Applying the ASSL approach to different human activity datasets. Besides, the ASSL
51 approach is compared with other well-known sequence learning approaches and the
52 results are presented.

53 The remainder of this paper is organised as follows: Section 2 discusses works related
54 to this paper. Section 3 describes the research methodology explaining an overview of the
55 proposed framework. In Section 4, the method proposed in this work for unsupervised
56 human activity segmentation is presented and Section 5 follows with a description of the
57 sequence learning method used in learning the activity segments constructed. Section 6
58 describes the application of the proposed model to human activity datasets and the results
59 obtained. In Section 7, the performance of the proposed ASSL is compared with other
60 sequence learning approaches and conclusions of the work are drawn in Section 8.

61 **2. Related Work**

62 There is a growing interest in research related to learning human activity sequences.
63 This section presents a review of relevant works in two categories: the segmentation of
64 human activities for detecting constituent actions, and activity modelling through sequential
65 learning/prediction.

66 *2.1. Action Detection and Segmentation*

67 Recently, human activity recognition has received much attention with a lot of research
68 undertaken for its applications in different areas (Adama et al., 2018; Lara and Labrador,
69 2013; Presti and Cascia, 2016). Most of the proposed activity recognition models (Presti and
70 Cascia, 2016) can attain impressive performances in their respective areas of application.
71 The majority focus on supervised approaches to activity recognition in which there is a
72 sufficient amount of labelled data available to build training models. However, in real-
73 world situations where obtaining labels for activities is a rather daunting task, supervised
74 methods for activity recognition may not be feasible (Adama et al., 2019). On the other

75 hand, unsupervised learning methods, like clustering (Comaniciu and Meer, 2002) are best
76 suited for such applications.

77 An aspect of activity recognition which tends to be a challenge for many systems is
78 detecting underlying/constituents actions in activities. This information is important in
79 determining the structure of activities which is important when considering trends or
80 sequences in such activities (Li and Fu, 2014). Therefore, segmentation is performed on
81 data to obtain partitions which represent certain characteristics in activities. This is a
82 vital step in investigating activity sequences. Existing approaches to segmentation of
83 human activity differ in terms of the following categories (Aminikhanghahi and Cook,
84 2017; Aminikhanghahi and Cook, 2019); the activity types that are modelled, the sensing
85 technology used to acquire information and the computational intelligence methods used
86 in the segmentation process.

87 With a focus on human activity recognition from 3D human skeleton joints information,
88 i.e. the joint positions or angles, different methods have been proposed for detecting actions
89 in an activity. The authors in (Li and Fu, 2014) proposed a method for detecting atomic
90 actions which they call *actionlets* using motion velocity. The method combined the Harris
91 corner detector and Lucas Kanade (LK) optical flow to get velocity magnitudes. In our
92 previous work (Adama et al., 2019), a key frame extraction method using the combined
93 motion energy of all body joints in an activity has been proposed. Other works using the
94 kinetic energy poses to determine key poses in activities are found in (Nunes et al., 2017;
95 Shan and Akella, 2014). These methods then apply different machine learning algorithms
96 for classification of actions obtained for activity recognition.

97 *2.2. Sequential Modelling of Activities*

98 The study of sequence learning algorithms are reported (Suresh et al., 2010; Cui et al.,
99 2016; Wen and Wang, 2017; Zhu et al., 2018). Sequence learning algorithms are used for
100 the analysis of patterns generated through a series of observed information for recognition
101 or classification of activities (Zhu et al., 2018). Machine learning researchers have studied
102 sequence learning over so many decades. This led to the development of statistical models
103 such as Hidden Markov Models (HMM) (Fine et al., 1998; Rabiner and Juang, 1986) and
104 Autoregressive Integrated Moving Average (ARIMA) (Durbin and Koopman, 2012) which
105 were introduced for time series and temporal pattern recognition problems (Cui et al., 2016).
106 Recurrent Neural Networks (RNNs) have since evolved to solve sequence prediction problems

107 due to their recurrent lateral structure. Long-Short Term Memory (LSTM), a type of RNN,
108 have a unique ability to selectively pass information across time and are able to model
109 significantly long-term dependencies due to the gating mechanism they possess (Hochreiter
110 and Schmidhuber, 1997). LSTMs also can deal with the vanishing gradient problem. This
111 has seen impressive performances in a variety of real-world applications.

112 Concerning human activities, attempts to model human activity sequences have been
113 studied by various researchers (Wen and Wang, 2017; Medina-Quero et al., 2018) using
114 different temporal models for human activities recognition. HMM is used over predefined
115 motion features of 3D joint positions to learn the dynamics of human actions (Lv and
116 Nevatia, 2006). Conditional Random Field (CRF) is another generative model employed
117 in modelling human actions. The CRF is used in (Han et al., 2010) to estimate motion
118 patterns that correspond to manifold subspace of 3D joint position features for human
119 action recognition. Similar approaches employing generative models to model activities are
120 also proposed in (Shan and Akella, 2014; Ofli et al., 2014). The 3D joint positions obtained
121 through skeleton tracking tend to be noisy. Therefore, when the change between actions is
122 small, without the accurate selection of features, recognising precise action states becomes
123 difficult. This tends to undermine the performance of generative models. Such models
124 require an adequate amount of data for training as they are prone to over-fitting. Dynamic
125 Time Warping (DTW) (Choi and Kim, 2018) is another solution used in modelling actions
126 by defining the distance between two temporal sequences of activity actions. The learning
127 can then be achieved through nearest-neighbour classification. However, the performance
128 of DTW is dependent on a good measure of the samples similarity. It could also suffer from
129 temporal misalignment when handling periodic actions which could lead to degrading its
130 performance (Li and Prakash, 2011). Reyes-Ortiz et al. (2016) have proposed a Transition-
131 Aware Human Activity Recognition (TAHAR) system for the real-time classification of
132 physical human activities. The system combined the probabilistic output of consecutive
133 activity predictions of a Support Vector Machine (SVM) with a heuristic filtering approach
134 to address issues regarding the occurrence of transitions between activities and unknown
135 activities to the proposed learning algorithm. From their results, the system was able to
136 situations with and without activities transition information. Similar works for sequential
137 learning of human activities employing LSTM RNN are seen in the works by Liu et al.
138 (2016) and Li et al. (2017).

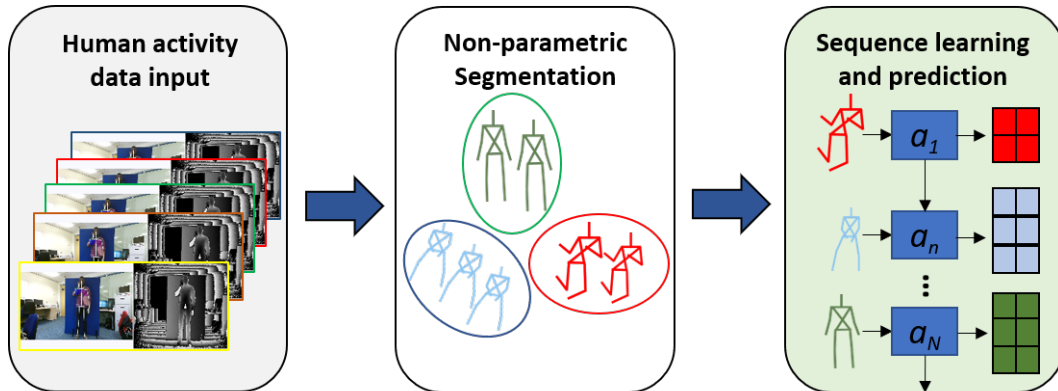


Fig. 2. Overview of the proposed approach to the Adaptive Segmentation and Sequence Learning (ASSL) of human activity.

139 These works demonstrate the effectiveness of segmentation and sequence modelling in
 140 exploring the underlying patterns in sequential data. This paper extends the approach
 141 of detecting key actions proposed in (Adama et al., 2019) for the segmentation of human
 142 activities by proposing an ASSL approach. Following from the identification of key actions,
 143 the non-parametric segmentation of 3D skeletal data of human activities obtained. This is
 144 then used in an LSTM model for the prediction of activity actions. In the following section,
 145 the problem statement is described and key definitions used in this work are presented.

146 3. Methodology

147 To address the challenges of segmentation and sequence learning of human activities,
 148 a novel framework for Adaptive Segmentation and Sequence Learning (ASSL) is proposed
 149 using visual information of activities. An overview of the ASSL framework is depicted in
 150 Fig. 2. There are three distinct steps in the proposed ASSL framework as described below:

- 151 1. Initially, key actions from observed human activity information are obtained. Human
 152 activities contain a large number of actions for which only the key aspects are relevant.
 153 By exploiting the temporal accumulated motion energy of each action through the
 154 sequence, the key actions can be extracted from the points of change in acceleration
 155 and deceleration of activity motion.
- 156 2. While segments of activities can be inferred from manual annotations, this creates a
 157 burden in *supervised* situations where high-dimensional data would require large

158 amounts of annotations to obtain feasible segments which can be learned. A
 159 non-parametric technique for feature space analysis is applied for *unsupervised*
 160 segmentation of relevant activity actions.

161 3. From the segments obtained, a Recurrent Neural Network (RNN) method for sequence
 162 learning called Long Short-Term Memory (LSTM) is used to learn activity sequences.

163 This work will benefit expert systems applications which require learning the underlying
 164 sequences in human actions through activities.

165 3.1. Definitions

166 Given a set of observed human activities $A = \{a_1, a_2, \dots, a_n, \dots, a_N\}$ performed by
 167 actors. The observations are obtained using an RGB-depth sensor. Each demonstration of
 168 an activity a_n within the observed activities set is a discrete time sequence of activity poses.
 169 An activity pose J as represented by;

$$J = [j_1, j_2, \dots, j_m, \dots, j_M], \quad \text{for } J \in \mathbb{R}^{3 \times M}, \quad (1)$$

170 is a feature space which represents 3D human skeleton joints with coordinates. M
 171 represents the total number of joints in J with each joint, j_m , with coordinates x_m, y_m, z_m
 172 corresponding to horizontal, vertical and depth positions respectively.

173 **Definition 1.** Key action, \bar{J} is defined as the important atomic level action performed
 174 during an activity. Key actions extracted from an activity represent a subset of poses
 175 $\bar{J} \subset a_n$, for $n = 1, 2, \dots, N$, which occurs in varying time instants of an executed activity.

176 **Definition 2.** Activity segmentation is defined by a function C in which each key action,
 177 \bar{J}_b , $b = 1, 2, \dots, B$, of an activity a_n is assigned a value, Q_z , $z = 1, 2, \dots, Z$, corresponding to
 178 a unique activity segment represented as:

$$C : a_n \mapsto (\bar{J}_b)_{1,2,\dots,B}, \quad \text{for } \bar{J}_b \in Q_z \quad (2)$$

179 where b is the index of the key action through the activity sequence and B is the number
 180 of key actions contained in a_n . Each segment derived comprises similar activity key actions
 181 through a temporal sequence.

182 **Definition 3.** Activity action sequence, S , is defined as the temporal ordering of all B key
 183 actions obtained from activity a_n . A repetition of similar key actions may be observed in
 184 the sequence at points where a_n contains actions which are repeated at different temporal
 185 instances. A representation of this definition is presented as:

$$S = (\bar{J}_b)_{b=1}^B \quad (3)$$

186 3.2. Assumptions

187 For the research presented in this paper, certain assumptions are made. They are:

- 188 - The observed sequence of a human activity comprises of unlabelled atomic actions
 189 which this work aims to identify through a process of adaptive segmentation.
- 190 - The number of key poses \bar{J}_B that make up an activity is not given. This is drawn
 191 from the fact that each activity can be segmented into key poses which make up for
 192 the relevant aspects that define the activity. However, this number is not pre-defined
 193 from activity observations in the proposed model.

194 3.3. Problem Statement 1

195 Given an observed sequence of human activity obtained using an RGB-depth sensor, the
 196 first phase is the segmentation of an unlabelled sequence into meaningful representations
 197 of similar actionlets. The segments obtained represent a collection of similar actions which
 198 may (or may not) fulfil temporal order relationship constraints.

199 The task of segmentation from an unlabelled activity sequence is addressed in this work
 200 using an adaptive approach to segmentation. The following steps are proposed for use in
 201 obtaining the function C for the segmentation of an activity.

202 **Detection of key actions (or poses):** Key actions of an activity are relevant in the
 203 process of learning an activity sequence. This is mainly because an activity can be executed
 204 in different forms whilst certain key aspects through the observation of an activity can
 205 uniquely identify it. As mentioned in the Introduction section, the motion energy feature
 206 of actions through an activity can be used in obtaining these key actions. The key actions
 207 are therefore identified by applying a filtering method of moving average crossovers of the
 208 motion energy. The description of how this is implemented is presented in the next Section.

209 **Non-parametric feature space clustering:** The key actions obtained from the filtering
210 process of the motion energy feature are clustered using a Mean-Shift feature space analysis
211 method. This method performs the clustering in terms of similarity of the motion energy
212 of key actions.

213 *3.4. Problem Statement 2*

214 To learn the sequence S of transition of actions from one activity segment to another,
215 it is important to note that an activity is not executed in only one possible sequence. An
216 activity can be executed with different temporal orders of constituent actions. This results
217 in a challenge of learning a generalised sequence for an activity.

218 The sequence of actions from one segment to another occur in intervals. The LSTM-
219 RNN algorithm, which is predominantly used in predicting time series, is applied in learning
220 the sequence of distinct actions within the activity segments. This method is used as the
221 algorithm is able to capture infinitely long sequences and predict succeeding occurrences
222 based on the memory gates.

223

224

225 The architecture of the ASSL approach for human activities from 3D skeleton information
226 as proposed in this paper is depicted in Fig. 3. This comprises three stages of activity data
227 input from an RGB-Depth sensor, segmentation of human activity and sequential learning
228 and prediction. Details of these stages are provided in the proceeding sections.

229 **4. Activity Segmentation**

230 Segmentation of human activity is relevant in the analysis of trends in transitions from
231 one activity state to another. This section describes the process of activity segmentation
232 using the extracted human activity information.

233 *4.1. Key Action Point Detection with Motion Energy*

234 Human activity consists of movement sequences generated by different body parts. It is
235 worth noting that not all aspects of an activity movement sequence are necessary to define
236 an activity. Certain aspects of the sequence can be executed in different forms and still result
237 in a similar activity. To simplify an activity to the relevant action points that constitute the

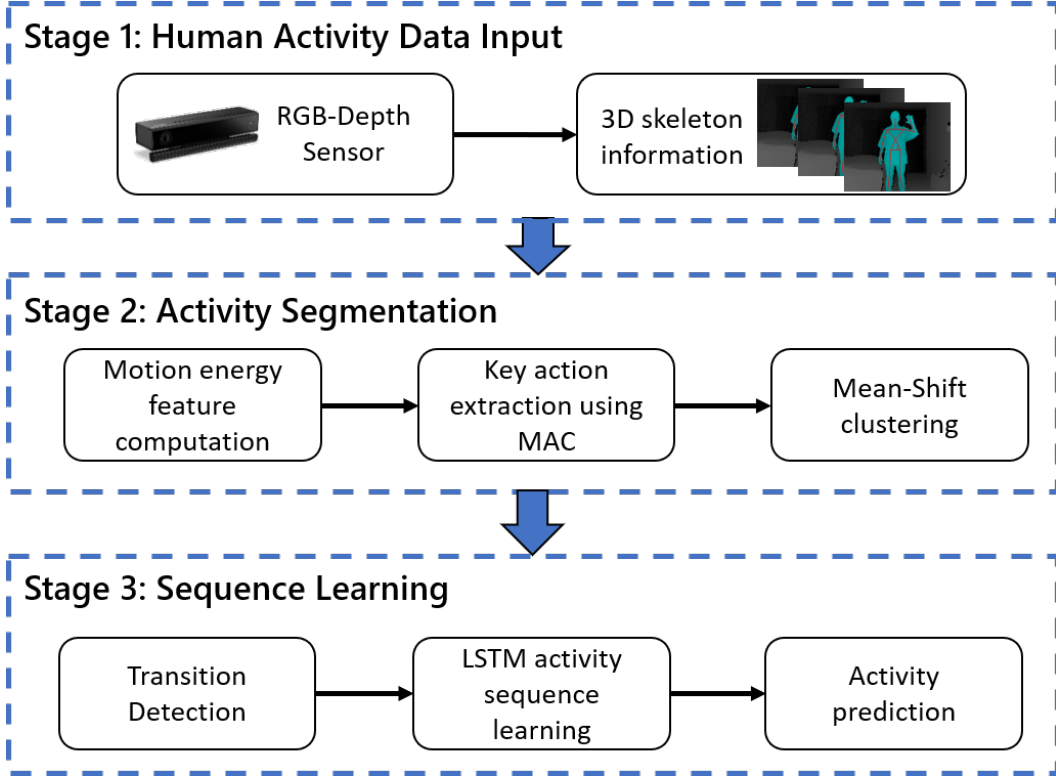


Fig. 3. Architecture of the proposed ASSL approach for human activities from 3D skeleton information which comprises activity input, segmentation and sequence learning stages respectively.

238 sequence, key poses are selected. This is achieved by leveraging the motion energy obtained
 239 from activity sequences.

240 4.1.1. Extraction of Motion Energy

241 The motion energy of activity poses as first proposed by (Shan and Akella, 2014) is based
 242 on the fact that joints show changes in acceleration and deceleration through an activity.
 243 This information is significant when considering the identification of the key action points
 244 of activities. Following from the representation of an activity pose given in Equation 1, the
 245 motion energy E_l for each pose is computed as the sum of motion energies for each joint in
 246 the pose;

$$E_l(J) = \sum_{m=1}^M E_l(j_m) \quad (4)$$

247 where j_m is a joint in the pose. It is assumed that the mass of all joints to be equally

248 one unit due to the fact that it is impossible to obtain the actual mass of a joint from the
 249 information obtained using RGB-Depth sensors. Computing the joint velocities using the
 250 temporal change ΔT in the position d of joints during an activity, the motion energy can
 251 be expressed as:

$$E_l(J) = \frac{1}{2} \sum_{m=1}^M (v_{j_m})^2 \quad (5)$$

252 where, v_{j_m} represents the velocity of joint j_m and is expressed as $v_{j_m} = \frac{d_m^c - d_m^p}{\Delta T}$, d_m^c is the
 253 current joint position and d_m^p is the previous joint position. By substituting v_{j_m} in Equation
 254 5, the motion energy of each joint is computed using the following equation:

$$E_l(J) = \frac{1}{2} \sum_{m=1}^M \left(\frac{d_m^c - d_m^p}{\Delta T} \right)^2 \quad (6)$$

255 4.1.2. Moving Average Crossover of Motion Energy

256 The Moving Average (MA) is a filtering technique often applied to get overall trends in
 257 data. This technique is used to highlight long-term cycles in time series data by smoothing
 258 out short-term variations (Droke, 2001). It works by creating series of averages of different
 259 time windows from a dataset over a given distribution.

260 Most of the works employing motion energy for identifying key action points of activities
 261 set threshold values of energy from a random exploration of selected points in order to
 262 extract the relevant points of interest in an activity (Nunes et al., 2017; Shan and Akella,
 263 2014; Zhu et al., 2015). The energy thresholds are selected by repeated experiments of
 264 different threshold values and the observations below the threshold value are selected as key
 265 poses. The MA of the extracted motion energy of poses are used in filtering the motion
 266 energy signal extracted from an activity sequence.

267 A different approach is proposed to use crossovers of two Simple Moving Averages
 268 (SMAs) of the extracted motion energy in identifying the relevant key poses of an activity.
 269 The SMA is an un-weighted mean of a set of data points. This is taken from equal sets of
 270 data to ensure variations in the mean and data points are aligned and not shifted in time.
 271 Given the motion energy obtained in Equation 4, the SMA for the motion energy signal of
 272 an activity can be computed using the following expression:

$$SMA = \frac{\sum_{r=0}^{\alpha-1} E_l(J)_{t-r}}{\alpha} \quad (7)$$

273 where α is the value of the period selected for MA and $t - r$ is the position of the selected
 274 observation within α . This is expressed in a simplified form as follows;

$$SMA_{E_l} = \frac{E_l(J)_t + E_l(J)_{t-1} + \dots + E_l(J)_{t-(\alpha-1)}}{\alpha} \quad (8)$$

275 Two moving averages are selected in this work - a short-term average (fast moving
 276 average) α_f and a long-term moving average (slow moving average) α_s . The MA crossovers
 277 are obtained from points where the SMAs for both α_f and α_s intersect. These points indicate
 278 significant changes in motion energy of activity poses and are used as reference points
 279 for their corresponding key actions in an activity sequence as presented in the following
 280 equation.

$$\overline{J}_b = SMA_{\alpha_s} \cap SMA_{\alpha_f} \quad (9)$$

281 Following the acquisition of the key action points, activity segments are obtained by
 282 application of a non-parametric feature space analysis technique - In this case, mean-shift
 283 clustering for associating key actions to clusters of similar actions.

284 4.2. Non-Parametric Clustering for Segmentation

285 Prior to learning the sequence of actions in an activity for prediction, it is necessary to
 286 know the segments that make up an activity. This information is not easily determined by
 287 mere observation of the key actions obtained from exploration of the motion energy feature.
 288 Also, the number of segments is defined for an activity as these can vary depending on
 289 the sequence observed. Therefore, the use of a non-parametric method of clustering key
 290 actions is proposed to determine the number of segments in an activity sequence and assign
 291 the obtained key actions to their respective segments before learning can be achieved. A
 292 mean-shift clustering approach is adopted here (Comaniciu and Meer, 2002). The mean-shift
 293 approach builds upon the concept of Kernel Density Estimation (KDE) (Parzen, 1962) which
 294 estimates the hidden distribution for a dataset by placing a kernel on each point contained
 295 in the dataset. The description of the mode of application for the proposed segmentation
 296 of human activity is provided below.

297 Given B key action points, \overline{J}_b $b = 1, \dots, B$, on a 2-dimensional space computed for an
 298 activity. As described in Section 4.1, these points correspond to the motion energies of key

Algorithm 1 Segmentation of human activity from joints coordinate skeleton information.

Input:

Instances of 3D skeleton joints coordinate of human activities $A = \{a_1, a_2, \dots, a_N\}$, in which each observation of activity a is a pose $J = [j_1, j_2, \dots, j_M]$;
 Activity time window t ;
 Moving average periods α_s, α_f ;

Output:

Activity segments obtained as a function C for assigning each key action to a segment;

Procedure:

- 1: **for** $a = 1$ to N **do**
 - 2: Find the velocity of each observation J within t ;
 - 3: Compute the motion energy for J : $E_l(J) = \sum_{m=1}^M E_l(j_m)$;
 - 4: Compute the simple moving average of the motion energy with the periods α_s, α_f :
 $SMA = \frac{\sum_{r=0}^{\alpha-1} E_l(J)_{t-r}}{\alpha}$;
 - 5: Key action points, $\bar{J}_b = SMA_{\alpha_s} \cap SMA_{\alpha_f}$;
 - 6: **end for**
 - 7: Assign \bar{J}_b to a cluster Q_z which is determined by a non-parametric mean-shift clustering technique;
 - 8: **return** $Q_Z = C(\bar{J}_b)$.
-

299 action positions. The kernel density estimate for the key action points with kernel K with
 300 a bandwidth parameter h is:

$$f(\bar{J}) = \frac{1}{Bh^2} \sum_{b=1}^B K\left(\frac{\bar{J} - \bar{J}_b}{h}\right) \quad (10)$$

301 with K satisfying the following two conditions:

- 302 1. $\int K(\bar{J})d\bar{J} = 1$, and
- 303 2. $K(\bar{J}) = K(|\bar{J}|)$ for all values of \bar{J} .

304 The first condition is required to ensure normalisation of the density estimate while the
 305 second condition relates to the symmetry of the data space containing all key action points
 306 of an activity. By applying a Gaussian symmetric kernel function for $K(\bar{J})$, the gradient of
 307 the density estimator in Equation 10 takes the form:

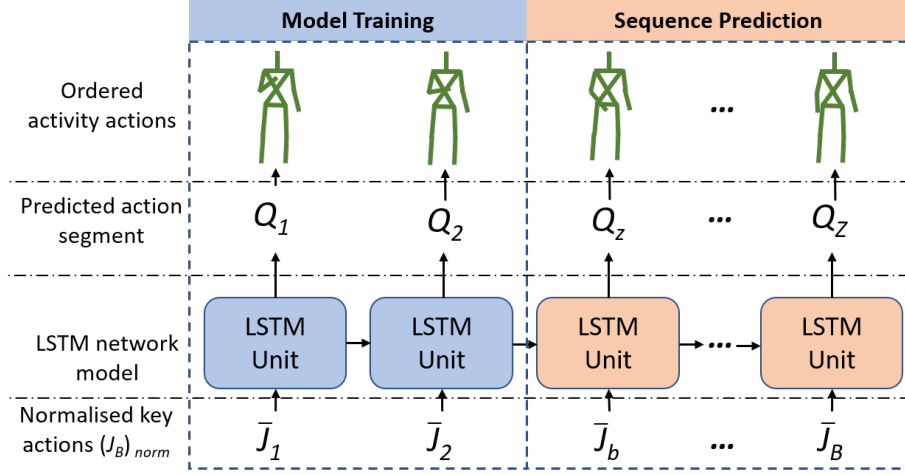


Fig. 4. LSTM structure for sequential learning and prediction of key action segments of human activity.

$$\nabla f(\bar{J}) = \frac{2}{Bh^4} \left(\sum_{b=1}^B g \left(\left| \frac{\bar{J} - \bar{J}_b}{h} \right| \right) \right) \vec{X}(\bar{J}) \quad (11)$$

where $\vec{X}(\bar{J})$ is the mean-shift vector pointing in the direction of increasing density and is represented as:

$$\vec{X}(\bar{J}) = \left(\frac{\sum_{b=1}^B \bar{J}_b g \left(\left| \frac{\bar{J} - \bar{J}_b}{h} \right| \right)}{\sum_{b=1}^B g \left(\left| \frac{\bar{J} - \bar{J}_b}{h} \right| \right)} - \bar{J} \right) \quad (12)$$

and $g(|\bar{J}|)$ is the derivative of the Gaussian kernel.

With the KDE computed, the mean-shift procedure is carried out by successive;

- Computation of the mean-shift vector $\vec{X}(\bar{J}_b)$ at the location of each key action point \bar{J}_b ,
- Translation of each action point $\bar{J}_b \rightarrow \bar{J}_b + \vec{X}(\bar{J}_b)$,
- Repeat until convergence, that is, where the gradient density function is zero.

Afterwards, the key action points identified at the same points are segmented as belonging to the same cluster Q_z . For further details of convergence, readers are referred to (Comaniciu and Meer, 2002). Algorithm 1 list the procedure for activity segmentation proposed in this paper.

320 **5. Sequence Learning and Prediction Model**

321 The sequence learning stage involves the learning of activity sequences from the
 322 segmented key actions. An LSTM network (Hochreiter and Schmidhuber, 1997) is used to
 323 learn the long-term contextual dependencies between key actions of an activity. The
 324 segmented key actions are used as input to the network for learning the dependencies
 325 between the action segments. This is further extended to predicting sequential actions of
 326 activities. Fig. 4 illustrates the structure of an LSTM network as applied in this work.
 327 The LSTM comprises of the following components: input gate i_t , forget gate f_t , a cell with
 328 a self-recurrent connection and output gate o_t . The key action segments obtained for an
 329 activity are normalised for standardisation of the values, thus resulting in
 330 $Q_{norm} = \{\bar{J}_{1Q_1}, \dots, \bar{J}_{BQ_z}\}_{norm}$. By taking Q_{norm} as input to the network, the network is
 331 updated every t timestep by iterating through all instances of the normalised key actions
 332 using the following equations;

$$i_t = \sigma(W^i(\bar{J}_{bQ_z}(t)) + U^i H_{t-1} + V^i) \quad (13)$$

$$f_t = \sigma(W^f(\bar{J}_{bQ_z}(t)) + U^f H_{t-1} + V^f) \quad (14)$$

$$o_t = \sigma(W^o(\bar{J}_{bQ_z}(t)) + U^o H_{t-1} + V^o) \quad (15)$$

$$g_t = \tan H(W^g(\bar{J}_{bQ_z}(t)) + U^g H_{t-1} + V^g) \quad (16)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (17)$$

$$H_t = o_t \odot \tan H(c_t) \quad (18)$$

333 where, $\sigma(\cdot)$ and $\tan H(\cdot)$ are the sigmoid and hyperbolic functions respectively. W, U, V are
 334 parameters of the LSTM model. The operation \odot denotes the element-wise multiplication
 335 of two vectors. The use of LSTM is due to its ability to map input activity sequences by
 336 recursively transforming current inputs Q_{norm} with the output hidden vector of previous
 337 steps H_{t-1} . Also, the vanish gradient problem inherent with RNN's is overcome by the
 338 memory cell c_t which is computed, allowing the error derivatives to flow in a different path.

339 **6. Application of Adaptive Segmentation and Sequence Learning Framework to** 340 **3D Skeleton Data of Daily Human Activity**

341 This section reports the experimental procedure and results of applications of the
342 proposed ASSL framework on 3D skeleton human activity datasets. To illustrate the
343 application of the proposed work of ASSL of human activity sequences, the model
344 proposed was applied to selected human activities. The proposed model is adaptive to
345 different activities and thus gives it the ability to deal with complexities in activities.

346 To understand the methodology and its ability to solve the problems identified in the
347 earlier Sections 3.3 and 3.4, the following hypotheses are proposed and evaluated.

348 **Hypothesis 1.** Where an unlabelled sequence of activity data is available, the segmentation
349 technique proposed can be used to identify unique segments of an activity used for label
350 assignments of actions in the sequence.

351 **Hypothesis 2.** Activity segments identified can be used to learn sequences for prediction
352 with a reliable performance.

353 To address these hypotheses, two activities are selected from two real world human activity
354 datasets; Dataset 1 - An experimental human activity dataset collected for this work and
355 Dataset 2 - A benchmark public dataset, Cornell Activity Dataset (CAD-60) (Sung et al.,
356 2011).

357 *6.1. Experimental Design and Datasets*

358 The motivation for the proposed ASSL framework is to address the issue of unlabelled
359 sequences of human activities, in such cases where there is no knowledge *a priori* of
360 constituent actions and their order, whilst there is the need to develop a system for
361 identifying the patterns of activities. The experimental design and datasets used in
362 evaluating the proposed framework are presented in this section.

363 *6.1.1. Dataset 1 - Experimental Human Activity Dataset*

364 The dataset generated to evaluate the proposed system in this work consists an activity
365 which involves a person picking up an object placed on a surface. A Microsoft Kinect
366 version 2 RGB-Depth sensor is used to acquire the 3D joint coordinate information of
367 person. This information is obtained at 30 fps. This activity is chosen due to the proposed

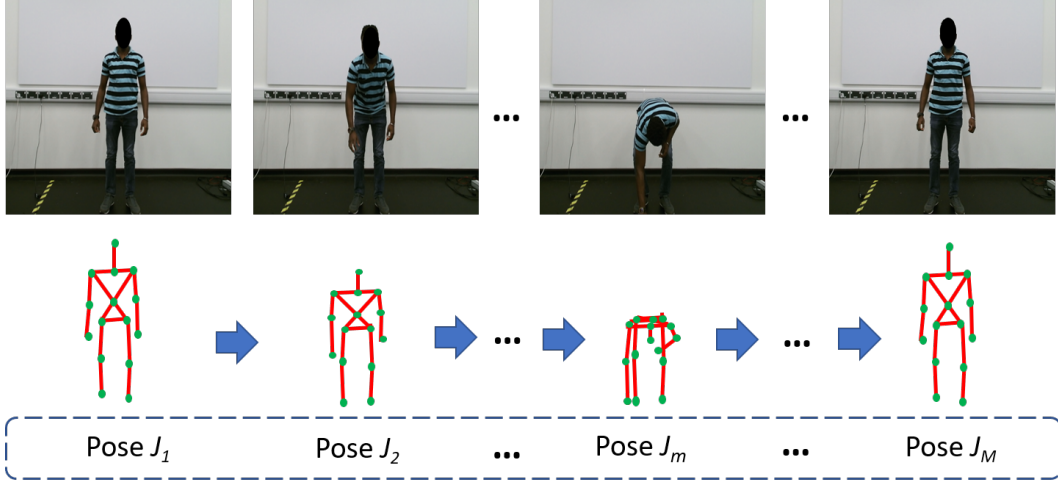


Fig. 5. Sample frames of *pick up object* activity obtained from the experimental activity dataset using an RGB-Depth sensor.

work being focused on enhancing the ability of assistive robots learning activity sequences for independent prediction of actions. Fig. 5 shows sample frames of the selected activity carried out by a person.

To obtain adequate amount of data to evaluate the ASSL framework, the activity is performed by three people. Each person is required to pick up an object from a flat surface repeatedly eight to ten times while the joint positions are recorded throughout the sequence. Table 1 shows the number of frames acquired from each person while carrying out the activity.

6.1.2. Dataset 2 - Cornell Activity Dataset (CAD-60)

The CAD-60 dataset (Sung et al., 2011) is based on human activity data obtained using an RGB-Depth sensor. The dataset comprises three modes of human activities data,

Table 1: Experimental dataset acquired from three actors for an activity - pick up object from a flat surface.

Activity	Number of frames			Total
	Person 1	Person 2	Person 3	
Pick up object	1804	1663	1355	4822

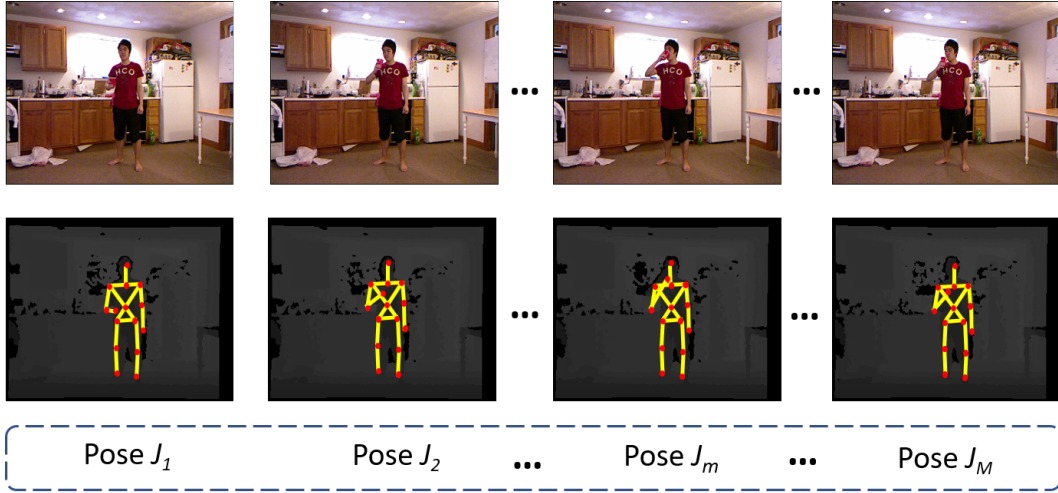


Fig. 6. Sample frames of *drinking water* activity obtained using an RGB-Depth sensor contained in the CAD-60 dataset (Sung et al., 2011). The sample shows RGB images and the corresponding depth image with the tracked skeleton overlaid.

379 RGB images, Depth images and 3D skeleton joint coordinates observed from a person
 380 performing an activity. The skeleton joint data consists of joint coordinates information of
 381 15 joints. The dataset is recorded at a frame rate of 15fps using a Microsoft Kinect sensor
 382 and includes recordings for 12 human activities namely; Rinsing mouth, brushing teeth,
 383 wearing contact lens, talking on the phone, drinking water, opening pill container, cooking
 384 (chopping), cooking (stirring), talking on couch, relaxing on couch, writing on whiteboard,
 385 working on computer and a sequence of random plus stationary activities. The data is
 386 collected from four participants with each performing each activity.

387 Most applications of this dataset are based on activity classification and therefore involve
 388 the use of all activities within the dataset. However, to demonstrate the work proposed in
 389 this paper, a single activity from the dataset is selected and used in our evaluations. The
 390 activity chosen is the *drinking water* activity as there are more motions involved in the
 391 activity when compared to the remainder activities available in the dataset. This creates
 392 a scenario with varying motion patterns to test the robustness of the framework. Sample
 393 frames of varying actions occurring throughout the activity sequence are shown in Fig. 6.
 394 The samples show a person drinking water with the tracked skeleton joints overlaid on the
 395 depth images. The activity is performed repeatedly 2 – 3 times.

396 *6.2. Experimental Human Activity Dataset Results and Evaluation*

397 To evaluate the performance of the proposed framework on the experimental dataset,
398 it is implemented in stages, starting with the segmentation process - the computation of
399 motion energy, detection of key action points and the non-parametric clustering for key
400 action segmentation. This is then followed by the sequence learning and prediction of the
401 obtained key actions.

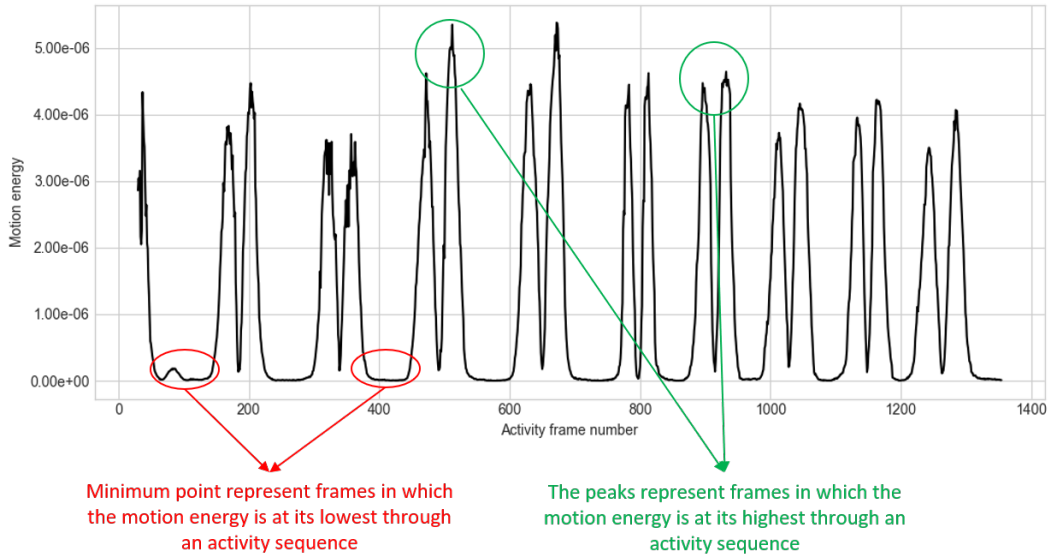
402 *6.2.1. Key Action Identification using Motion Energy*

403 Applying the approach to identifying key action points of an activity, the motion energy
404 is computed for 3D joint positions data obtained from each person. A window size, w_s , of
405 one second is used which corresponds to 30 frames of activity to compute the motion energy.
406 Fig. 7a shows the motion energy obtained from person 1 of the experimental dataset. The
407 figure shows the changes in the cumulative motion energy which is a result of continuous
408 acceleration and deceleration of body joints through the activity sequence.

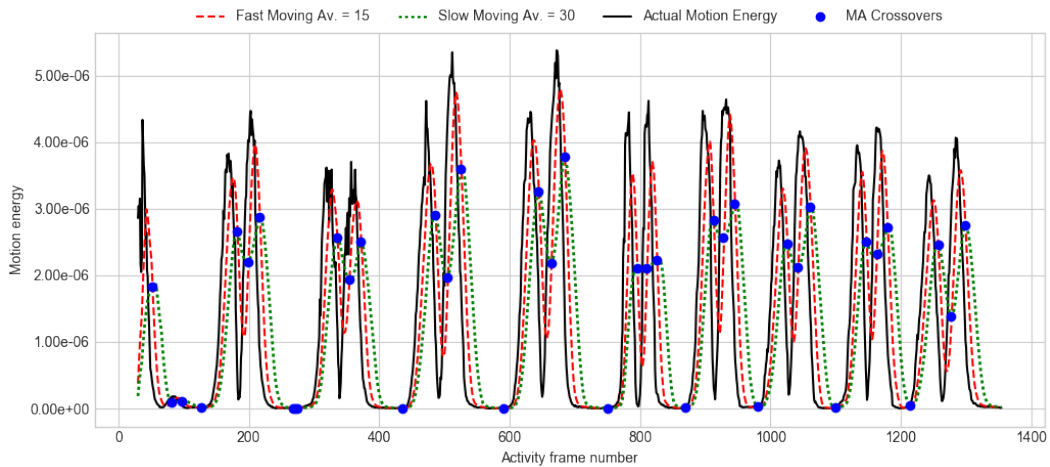
409 In the proposed framework, the key actions are identified at points of minimum and
410 maximum motion energies. Applying the simple moving average technique, after multiple
411 experiments with different values of SMA_{α_s} and SMA_{α_f} , 30 and 15 frames are selected for
412 both moving averages respectively. Fig. 7b depicts the key action points identified from the
413 motion energy computed in Fig. 7a. The green plot shows the SMA_{α_s} while the red plot
414 shows the SMA_{α_f} . The crossover points of both moving averages are identified by the blue
415 dots in Fig. 7b. These points represent the key actions $\overline{J_B}$ in the activity sequence from
416 the data. Similarly, the key actions are obtained for all participants in the experimental
417 dataset.

418 *6.2.2. Non-parametric Clustering of Experimental Dataset*

419 Due to the varying nature of the activities performed from one individual to another,
420 there are variations in motion energy values from person to person. To tackle this difficulty,
421 the motion energy of the key actions identified for each participant's activity are normalised
422 for standardisation across all participants. Fig. 8 shows the representation of normalised
423 motion energies of identified key actions for all persons in the dataset. A total of 202 key
424 action frames are identified from all three participants which shows a significant reduction
425 when compared to the total number of frames 4822 as shown in Table 1. This emphasises the



(a)



(b)

Fig. 7. Key action identification for *pick up object* activity from person 1 in the experimental dataset. (a) Motion energy plot for person 1 from the experimental dataset. The energy is computed using a 1 second window = 30 frames. (b) Motion energy plot with identified crossover points of two moving averages which represent the identified key action points of the activity. $SMA_{\alpha_f} = 15$ and $SMA_{\alpha_s} = 30$.

426 need for the segmentation process to reduce the computational complexities when learning
 427 the activity sequence.

428 The normalised values are then clustered using the non-parametric method described

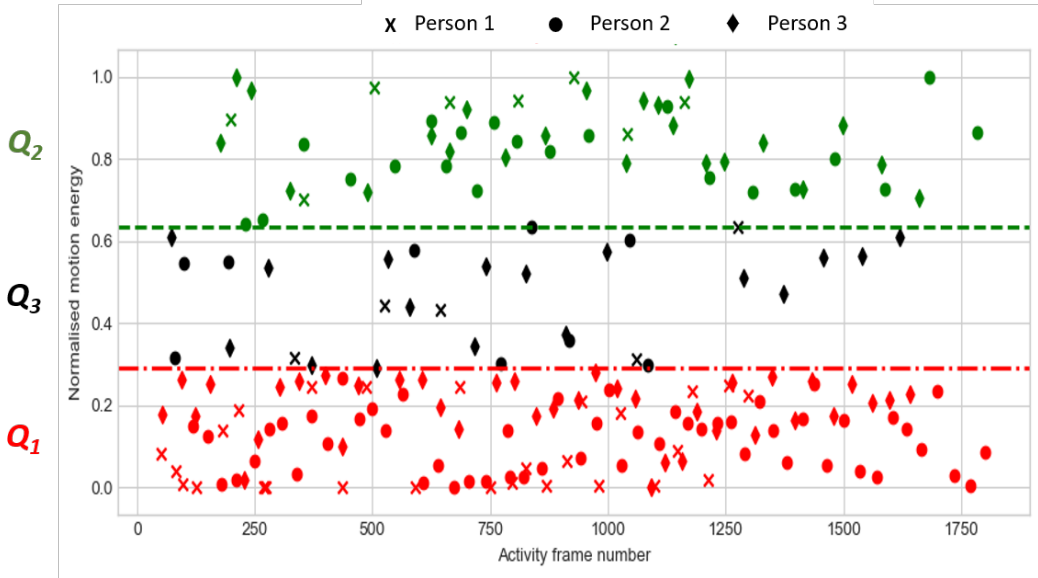


Fig. 8. Normalised motion energy with action segment identification of key actions for all participants in the experimental human activity dataset corresponding to the *pick up object activity*.

429 earlier. The results obtained from clustering is also represented in Fig. 8. It can be observed
 430 that for the selected activity three segments corresponding to Q_1 , Q_2 and Q_3 , are identified
 431 and the boundaries of the segments as obtained from the results are represented by the
 432 horizontal line plots (green and orange) shown on the figure. Fig. 9 shows the distribution
 433 of the number of key action points identified in each activity segment for all participants.

434 6.2.3. Sequence Learning of Experimental Human Activity Dataset

435 The sequence learning model is grounded on the results obtained from the activity
 436 segmentation process. To investigate the performance, the outputs from the segmentation
 437 process are fed as input to the learning model and a comparison is made between the
 438 results obtained and the actual activity sequence observed. This comparison is done in
 439 terms of the Mean Absolute Error (MAE), Mean Absolute Scaled Error (MASE) and
 440 Root-Mean-Square Error (RMSE) for the predictions made. The performance of the
 441 sequence learning model in this work depends on a proper segmentation of the unlabelled
 442 activity sequences observed.

443 The performance of the sequence learning framework is evaluated on the experimental
 444 dataset. Considering the dataset consists of 3 participants, a leave-one-out cross validation

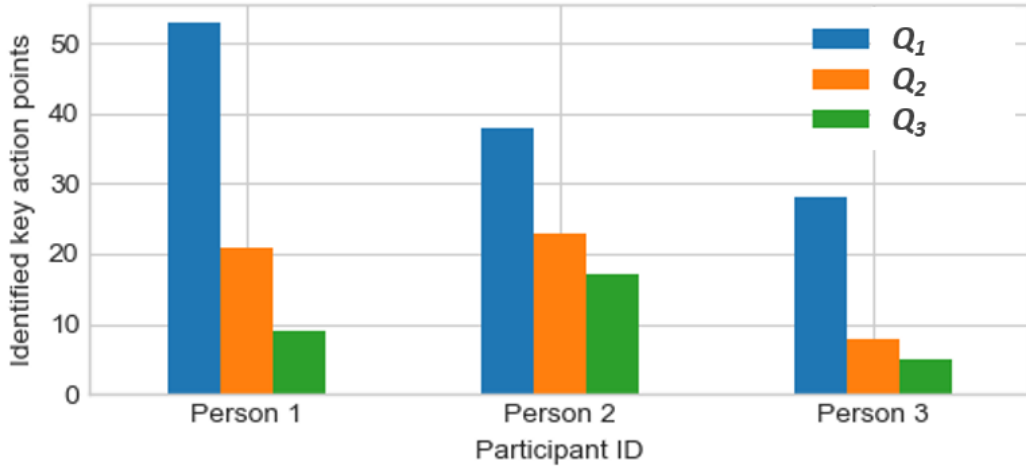


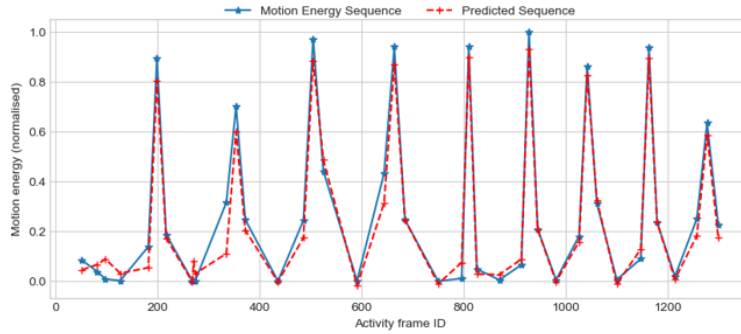
Fig. 9. Activity segmentation distribution for participants in the experimental human activity dataset.

445 approach is used in experiments to learn sequences of key action occurrences for an activity.
 446 Two participants are used in training the model and the remainder is left out for testing.
 447 This is done through consecutive iterations with each participant used in testing the model.

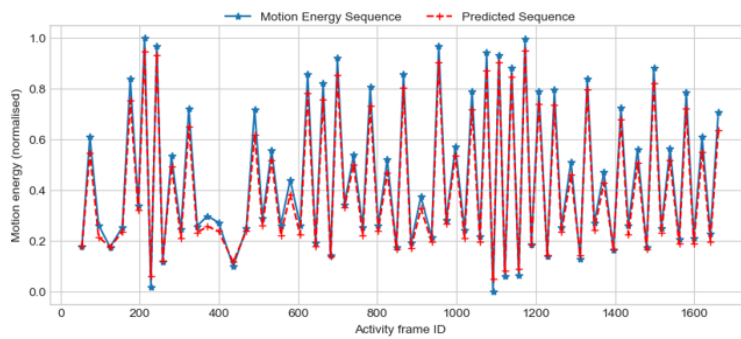
448 Fig. 10 shows the result of the sequence learning model on the prediction of the
 449 activity sequence contained in the experimental dataset. Table 2 shows the result when
 450 the experimental dataset is applied to the proposed ASSL model. The results produced
 451 RMSE values of 0.055, 0.049 and 0.050 respectively for all three participants in the dataset
 452 when each was tested using the leave-one-out cross validation. The lower the RMSE value
 453 the better the result in predicting the sequence. The variation in the structure of the
 454 sequence between the remainder two person’s data used when training the model and the
 455 structure of the person 1 used in testing the model produced a higher RMSE value (0.055)
 456 in comparison with the RMSE value obtained for the other two. This can be attributed to
 457 the nature of the activity sequence for person 1, i.e. the speed of the activity.

458 6.3. CAD-60 Dataset Results and Evaluation

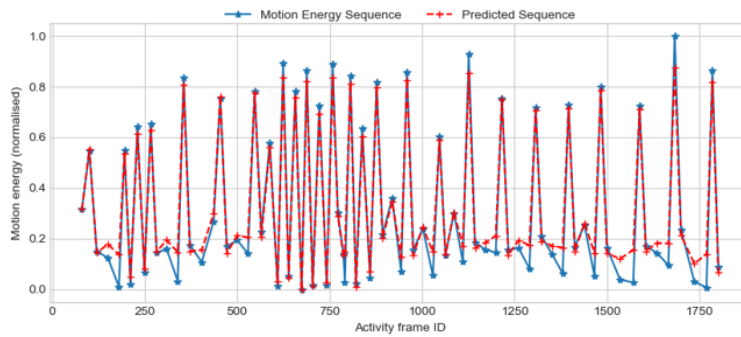
459 The segmentation process applied to the CAD-60 dataset using the same values of simple
 460 moving averages as in the case of the experimental activity dataset to identify key actions
 461 which are segmented resulted in a similar number of activity segments. The distribution
 462 of key actions identified in each segment is given Fig. 11. This shows a similar ratio in
 463 the distribution of key actions identified for all actors except for the case of *actor 1*. This



(a)



(b)



(c)

Fig. 10. Performance of sequence learning model on the prediction of experimental dataset activity sequence. (a) Person 1. (b) Person 2. (c) Person 3.

464 infers that for the activity - *drinking water* - performed by all actors, there are three atomic
 465 actions that define the activity. The order in which the actions occur define the activity
 466 sequence. It is important to note that the segments identified in the experiments with the
 467 CAD-60 experiment are not the same as those of the experimental activity dataset.

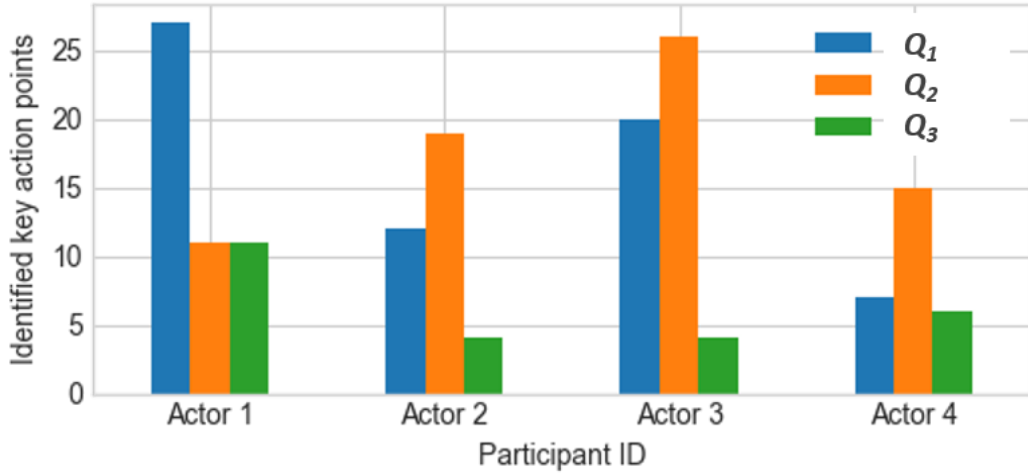


Fig. 11. Distribution of key action points in identified activity segments for all actors in the CAD-60 dataset.

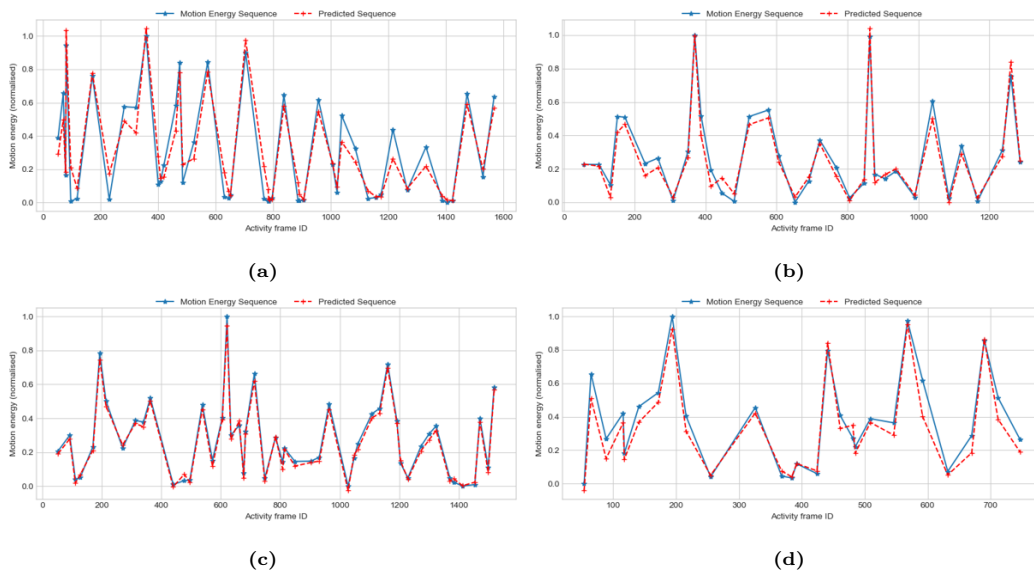


Fig. 12. Prediction performance of sequence learning model on the CAD-60 dataset. (a) Actor 1. (b) Actor 2. (c) Actor 3. (d) Actor 4.

468 Evaluating the performance of the sequence learning framework with the CAD-60 dataset
 469 is implemented in a similar method to the experimental dataset. A leave-one-out cross
 470 validation approach is also applied with each participant data used in testing while the
 471 remainder three are used in training the model. This is performed in consecutive iterations.

472 In Fig. 12, the prediction results for all actors are shown. The plots in the figure represent
473 when each actors' activity data is left out from the training process and used to test the
474 trained sequence learning model. Table 3 shows the prediction results obtained for the
475 dataset with the ASSL. The RMSE values produced from predicting activity sequences
476 for the data tested correspond to 0.092, 0.053, 0.025 and 0.078 for Actor's 1, 2, 3 and 4
477 respectively. The low RMSE values show the model is able to learn with a high degree of
478 reliability the activity sequence.

479 It should be noted that in the experiments a consideration was given to test the sequence
480 learning model without the process of segmentation to extract key actions, that is, using
481 the motion energy of all actions within the activity. This was done in the evaluation of
482 the proposed model. Using all the actions, the clustering stage identified the actions as
483 belonging to one cluster as opposed to the output of the clustering using the segmented key
484 actions. From the visual observation of the activity, it is clear that this activity consists of
485 more than one distinguishable action. Also, for both datasets used, the sequence learning
486 model performed poorly in predicting the action sequences. This could be due to all actions
487 identified as being the same.

488 **7. Comparison with other Sequence Learning Models**

489 This section presents a comparison of the proposed ASSL framework's performance
490 with other statistical models widely used in learning sequences from time series data. The
491 adaptive segmentation and sequence learning of 3D skeleton data of human activities
492 framework primarily demonstrates that unlabelled actions and sequences of activities can
493 be modelled for accurate prediction of unseen actions. This is beneficial for applications
494 that require exploiting the underlying patterns to understand human tasks from visual
495 observations while they are executed. This was demonstrated in the previous sections. To
496 further emphasise the ability of the proposed framework to learn activity sequences, a
497 comparison is made with other methods of sequence learning used in forecasting
498 applications, an Autoregressive Integrated Moving Average (ARIMA), Support Vector
499 Regression (SVR) (Gascon-Moreno et al., 2012; Awad and Khanna, 2015) and Gaussian
500 Mixture Regression (GMR). The basis for selecting the ARIMA model is because it comes
501 from a well established area of computational intelligence. ARIMA models are also widely
502 used in analysis of temporal pattern recognition and time series prediction. The SVR and

503 GMR models are techniques mostly applied in batch learning problems for forecasting
 504 purposes. These models are applied to both the experimental dataset and CAD-60 dataset
 505 described in Sections 6.1.1 and 6.1.2 respectively, with the same validation technique
 506 already described.

507 Autoregressive Moving Average (ARMA) models are amongst the most widely used
 508 statistical algorithms for modelling and predicting time series information (Smith et al.,
 509 2018). A generalisation of this model is the Autoregressive Integrated Moving Average
 510 (ARIMA) which is applied in situations where there is evidence of non-stationarity in data.
 511 In such cases, a differencing step, d , corresponding to the *Integrated* part of the model is
 512 applied to remove non-stationarity points (Ümit Çavus Büyüksahina and Ertekin, 2019).
 513 Afterwards, the ARMA model is applied on the stationary data. The implementation of
 514 ARIMA in this work follows the method described in (Ümit Çavus Büyüksahina and Ertekin,
 515 2019). The Auto-Regressive, *AR*, component uses weighted linear combinations of previous
 516 values of the data sequence and performs a regression of the sequence against itself. Similarly,
 517 the Moving Average, *MA*, component attempts predicting a target using regression based

Table 2: Comparison of the proposed ASSL model’s performance with ARIMA, SVR and GMR models on the experimental human activity dataset (the best results across all models in bold text).

Metric	Method	Person 1 (error \pm var.)	Person 2 (error \pm var.)	Person 3 (error \pm var.)
MAE	ASSL	0.044 \pm 0.005	0.025 \pm 0.006	0.032 \pm 0.004
	ARIMA	0.228 \pm 0.032	0.135 \pm 0.036	0.132 \pm 0.069
	SVR	0.057 \pm 0.005	0.076 \pm 0.006	0.072 \pm 0.006
	GMR	0.345 \pm 0.090	0.407 \pm 0.090	0.309 \pm 0.077
MASE	ASSL	0.152 \pm 0.005	0.122 \pm 0.006	0.047 \pm 0.004
	ARIMA	0.586 \pm 0.032	0.272 \pm 0.036	0.291 \pm 0.069
	SVR	0.141 \pm 0.005	0.153 \pm 0.006	0.244 \pm 0.006
	GMR	0.849 \pm 0.090	0.823 \pm 0.090	1.046 \pm 0.077
RMSE	ASSL	0.055 \pm 0.005	0.049 \pm 0.006	0.050 \pm 0.004
	ARIMA	0.298 \pm 0.032	0.198 \pm 0.036	0.175 \pm 0.069
	SVR	0.075 \pm 0.005	0.088 \pm 0.006	0.081 \pm 0.006
	GMR	0.457 \pm 0.090	0.506 \pm 0.090	0.414 \pm 0.077

Table 3: Comparison of the proposed ASSL model’s performance with ARIMA, SVR and GMR models on the CAD-60 dataset (the best results are in bold text).

Metric	Method	Actor 1	Actor 2	Actor 3	Actor 4
		(error \pm var.)	(error \pm var.)	(error \pm var.)	(error \pm var.)
MAE	ASSL	0.072 \pm 0.023	0.044 \pm 0.015	0.023 \pm 0.012	0.062 \pm 0.018
	ARIMA	0.307 \pm 0.190	0.202 \pm 0.077	0.220 \pm 0.109	0.255 \pm 0.122
	SVR	0.123 \pm 0.023	0.100 \pm 0.014	0.089 \pm 0.010	0.100 \pm 0.017
	GMR	0.302 \pm 0.117	0.273 \pm 0.062	0.239 \pm 0.050	0.357 \pm 0.093
MASE	ASSL	0.281 \pm 0.023	0.336 \pm 0.015	0.442 \pm 0.012	0.253 \pm 0.018
	ARIMA	0.865 \pm 0.190	0.690 \pm 0.077	0.983 \pm 0.109	0.802 \pm 0.122
	SVR	0.312 \pm 0.023	0.385 \pm 0.014	0.452 \pm 0.010	0.341 \pm 0.017
	GMR	0.765 \pm 0.117	1.045 \pm 0.062	1.208 \pm 0.050	1.215 \pm 0.093
RMSE	ASSL	0.092 \pm 0.023	0.053 \pm 0.015	0.025 \pm 0.012	0.078 \pm 0.018
	ARIMA	0.339 \pm 0.190	0.267 \pm 0.077	0.264 \pm 0.109	0.326 \pm 0.122
	SVR	0.153 \pm 0.023	0.119 \pm 0.014	0.105 \pm 0.010	0.130 \pm 0.017
	GMR	0.456 \pm 0.117	0.370 \pm 0.062	0.326 \pm 0.050	0.469 \pm 0.093

518 on past forecast errors. The parameters of the ARIMA model corresponding to coefficients
519 of the orders of the model are d , p and q . p represents the number of time lags to consider.
520 When $p = 0$, the mode is reduced to a MA model of q order. Similarly, if $q = 0$, the
521 model becomes AR of p order. Details of the selection of the optimal parameters for the
522 ARIMA model used are beyond the scope of this work. Readers are referred to (Ümit Çavus
523 Büyüksahina and Ertekin, 2019) for more insight into ARIMA.

524 The SVR model as a supervised learning approach, has been applied as an effective tool
525 in real-value function estimation and is characterised by the use of kernels. The model is
526 trained by using a symmetrical loss function which penalises high and low misestimates
527 equally. This model is used in the evaluation process to validate the proposed ASSL models
528 performance. Implementations of the SVR and GMR models follow the methods in (Sung,
529 2004; Awad and Khanna, 2015).

530 *7.1. Evaluation of the Results of ARIMA, SVR and GMR Prediction Models on the*
531 *Experimental Dataset*

532 The normalised key action points of the motion energy extracted from the experimental
533 human activity are used as input to the ARIMA, SVR and GMR models as mentioned
534 earlier. The results shown in Table 2 present the performance of all the models on the
535 experimental dataset. As observed from the table, the proposed ASSL model had a better
536 performance in terms of the MAE and RMSE than all the other models when observed across
537 all the participants in the dataset. There is a significant difference in the MAE and RMSE
538 performance obtained with the ASSL method outperforming all the other models. Next
539 to the MAE and RMSE performance of the ASSL, the SVR model obtained comparable
540 performance. However, the SVR model did slightly better than the ASSL model in terms
541 of the MASE performance for person 1. As with most unsupervised learning structures,
542 the ARIMA is able to predict data sequences with only the targeted data. It can also be
543 noted from the results of Table 2 that the GMR model had the least performance across all
544 the participants when compared with the other models. The GMR algorithm is known to
545 be a fast learning model as it maximises only the likelihood. However, when it encounters
546 many points, estimating the covariance matrices tends to be difficult. Therefore, the model
547 diverges.

548 *7.2. Evaluation of the Results of ARIMA, SVR and GMR models on CAD-60 Dataset*

549 Table 3 shows a comparison of the results obtained for the performance of the ASSL
550 framework with the ARIMA, SVR and GMR models on the CAD-60 dataset. Similar
551 to the performance obtained with the experimental dataset, the proposed ASSL model
552 outperformed all the other models with lower error values across all four actors. Similarly,
553 the GMR was the worst-performing model on the CAD-60 dataset, except for the MASE
554 for actor 1 where the ARIMA model had the highest error value.

555 The ARIMA model works as a regression model and therefore does not require labelled
556 samples. However, the proposed approach is able to obtain labels through a non-parametric
557 approach which is used in the later stage of sequence learning. This gives the ASSL method
558 an edge over the ARIMA.

559 8. Conclusion and Future Work

560 In this paper, a novel adaptive technique for the segmentation and sequential learning
561 of human activities is presented. The goal is to enable the discovery unknown activity
562 patterns for prediction of future actions in an activity sequence, especially, for use in assistive
563 robotics. Due to the dynamic nature of human behaviour, there are uncertainties associated
564 with modelling actions performed in an activity. This work focused on proposing a model
565 capable of adapting to variations that exist in actions through activity sequences. The use
566 of 3D skeleton joint data obtained with RGB-Depth sensors makes it possible to acquire
567 representations of actions for learning such activities.

568 The motion energy of skeleton joints is used as a feature in the segmentation process.
569 This is due to changes in acceleration and deceleration observed in skeleton joints through
570 a continuous sequence of activities. This feature is used in identifying key actions in an
571 activity sequence from the moving average crossovers of the computed motion energy. This
572 steps acts as filter stage as not all actions of an activity are relevant in predicting the
573 activity sequence. We leverage the ability of LSTM model in learning activity sequences
574 for predicting future actions of activities based on previous instances. The results show
575 the performance of the LSTM sequence learning model is better than the unsupervised
576 sequence learning approaches. Furthermore, learning sequences of activity from unlabelled
577 activity structures are addressed. The segmentation approach used to identify labels from
578 the structures made it possible to solve the unsupervised learning problem with a supervised
579 technique of learning sequences.

580 Due to the challenges in this research area, the work presented in this paper has some
581 limitations which will be addressed in future work. The work presented in this paper can
582 be extended to include more subjects used in the experiments. This is needed due to the
583 variation that exist from person to person performing an activity. This will add robustness
584 to the sequence learning models. Furthermore, more research will be done on improving
585 the performance of the sequence learning and prediction models in order to reduce the
586 predictions errors in the results. Specifically, other variants of the LSTM RNN such as
587 Bidirectional-LSTM will be tested.

588 **References**

- 589 D. A. Adama, A. Lotfi, C. Langensiepen, K. Lee, P. Trindade, *Soft Computing* 22 (2018)
590 7027–7039.
- 591 D. O. Anderez, A. Lotfi, A. Pourabdollah, *Expert Systems with Applications* 140 (2020).
- 592 S. Suresh, K. Dong, H. Kim, *Neurocomputing* 73 (2010) 3012 – 3019.
- 593 R. Lioutikov, G. Neumann, G. Maeda, J. Peters, *The International Journal of Robotics*
594 *Research* 36 (2017) 879–894.
- 595 S. Krishnan, A. Garg, S. Patil, C. Lea, G. Hager, P. Abbeel, K. Goldberg, *The International*
596 *Journal of Robotics Research* 36 (2017) 1595–1618.
- 597 O. D. Lara, M. A. Labrador, *IEEE Communications Surveys Tutorials* 15 (2013) 1192–1209.
- 598 L. L. Presti, M. L. Cascia, *Pattern Recognition* 53 (2016) 130 – 147.
- 599 D. A. Adama, A. Lotfi, R. Ranson, P. Trindade, in: 2nd UK-RAS Conference on Embedded
600 Intelligence (UK-RAS19), pp. 60–63.
- 601 D. Comaniciu, P. Meer, *IEEE Transactions on Pattern Analysis and Machine Intelligence*
602 24 (2002) 603–619.
- 603 K. Li, Y. Fu, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36 (2014)
604 1644–1657.
- 605 S. Aminikhanghahi, D. J. Cook, in: 2017 IEEE International Conference on Pervasive
606 Computing and Communications Workshops (PerCom Workshops), pp. 262–267.
- 607 S. Aminikhanghahi, D. J. Cook, *Pervasive and Mobile Computing* 53 (2019) 75 – 89.
- 608 D. A. Adama, A. Lotfi, C. Langensiepen, in: *Advances in Computational Intelligence*
609 *Systems*, Springer International Publishing, 2019, pp. 303–311.
- 610 U. M. Nunes, D. R. Faria, P. Peixoto, *Pattern Recognition Letters* 99 (2017) 21 – 31.
- 611 J. Shan, S. Akella, in: 2014 IEEE International Workshop on Advanced Robotics and its
612 Social Impacts, pp. 69–75.

613 Y. Cui, S. Ahmad, J. Hawkins, *Neural Computation* 28 (2016).

614 J. Wen, Z. Wang, *Expert Systems with Applications* 74 (2017) 19 – 28.

615 H. Zhu, H. Chen, R. Brown, *Journal of Biomedical Informatics* 84 (2018) 148 – 158.

616 S. Fine, Y. Singer, N. Tishby, *Machine Learning* 32 (1998) 41–62.

617 L. Rabiner, B. Juang, *IEEE ASSP Magazine* 3 (1986) 4–16.

618 J. Durbin, S. Koopman, *Time series analysis by state space methods*, Oxford University
619 Press, New York, 2012.

620 S. Hochreiter, J. Schmidhuber, *Neural Comput.* 9 (1997) 1735–1780.

621 J. Medina-Quero, S. Zhang, C. Nugent, M. Espinilla, *Expert Systems with Applications* 114
622 (2018) 441 – 453.

623 F. Lv, R. Nevatia, in: *Computer Vision – ECCV 2006*, Springer Berlin Heidelberg, Berlin,
624 Heidelberg, 2006, pp. 359–372.

625 L. Han, X. Wu, W. Liang, G. Hou, Y. Jia, *Image and Vision Computing* 28 (2010) 836 –
626 849.

627 F. Ofli, R. Chaudhry, G. Kurillo, R. Vidal, R. Bajcsy, *Journal of Visual Communication*
628 and Image Representation 25 (2014) 24 – 38.

629 H.-R. Choi, T. Kim, *Mathematical Problems in Engineering* 2018 (2018).

630 L. Li, B. A. Prakash, in: *Proceedings of the 28th International Conference on International*
631 *Conference on Machine Learning, ICML’11*, pp. 185–192.

632 J.-L. Reyes-Ortiz, L. Oneto, A. SamÁ , X. Parra, D. Anguita, *Neurocomputing* 171 (2016)
633 754 – 767.

634 J. Liu, A. Shahroudy, D. Xu, G. Wang, in: *Computer Vision – ECCV 2016*, pp. 816–833.

635 K. Li, X. Zhao, J. Bian, M. Tan, in: *2017 IEEE International Conference on Mechatronics*
636 *and Automation (ICMA)*, pp. 1556–1561.

637 C. Droke, *Moving Averages Simplified*, Marketplace Books, 2001.

- 638 G. Zhu, L. Zhang, P. Shen, J. Song, L. Zhi, K. Yi, in: 2015 IEEE International Conference
639 on Robotics and Biomimetics (ROBIO), pp. 1209–1214.
- 640 E. Parzen, *The Annals of Mathematical Statistics* 33 (1962) 1065–1076.
- 641 J. Sung, C. Ponce, B. Selman, A. Saxena, in: *Proceedings of the 16th AAAI Conference on*
642 *Plan, Activity, and Intent Recognition, AAAIWS'11-16*, AAAI Press, 2011, pp. 47–55.
- 643 J. Gascon-Moreno, S. Salcedo-Sanz, E. Ortiz-Garcia, J. Acevedo-Rodriguez, J. A. Portilla-
644 Figueras, *Expert Systems with Applications* 39 (2012) 8220 – 8227.
- 645 M. Awad, R. Khanna, *Support Vector Regression*, Apress, Berkeley, CA, 2015, pp. 67–80.
- 646 E. M. Smith, J. Smith, P. Legg, S. Francis, in: *Advances in Computational Intelligence*
647 *Systems*, pp. 191–202.
- 648 Ümit Çavus Büyüksahina, S. Ertekin, *Neurocomputing* 361 (2019) 151 – 163.
- 649 H. G. Sung, *Gaussian Mixture Regression and Classification*, Ph.D. thesis, Rice University,
650 2004.