



# Silhouette-based Human Action Recognition using Sequences of Key Poses

Alexandros Andre Chaaaraoui<sup>a,\*</sup>, Pau Climent-Pérez<sup>a</sup>, Francisco  
Flórez-Revuelta<sup>b</sup>

<sup>a</sup>*Department of Computing Technology, University of Alicante, P.O. Box 99, E-03080,  
Alicante, Spain*

<sup>b</sup>*Faculty of Science, Engineering and Computing, Kingston University, Penrhyn Road,  
KT1 2EE, Kingston upon Thames, United Kingdom*

---

## Abstract

In this paper, a human action recognition method is presented in which pose representation is based on the contour points of the human silhouette and actions are learned by making use of sequences of multi-view key poses. Our contribution is two-fold. Firstly, our approach achieves state-of-the-art success rates without compromising the speed of the recognition process and therefore showing suitability for online recognition and real-time scenarios. Secondly, dissimilarities among different actors performing the same action are handled by taking into account variations in shape (shifting the test data to the known domain of key poses) and speed (considering inconsistent time scales in the classification). Experimental results on the publicly available

---

\*Corresponding author: Alexandros Andre Chaaaraoui, Department of Computing Technology, University of Alicante, P.O. Box 99, E-03080, Alicante, Spain. Phone: +34 965903681, Fax: +34 965909643

*Email addresses:* [alexandros@dtic.ua.es](mailto:alexandros@dtic.ua.es) (Alexandros Andre Chaaaraoui),  
[pcliment@dtic.ua.es](mailto:pcliment@dtic.ua.es) (Pau Climent-Pérez), [F.Florez@kingston.ac.uk](mailto:F.Florez@kingston.ac.uk) (Francisco  
Flórez-Revuelta)

*URL:* <http://www.dtic.ua.es> (Alexandros Andre Chaaaraoui)

Weizmann, MuHAVi and IXMAS datasets return high and stable success rates, achieving, to the best of our knowledge, the best rate so far on the MuHAVi *Novel Actor* test.

*Keywords:* human action recognition, key pose, key pose sequence, Weizmann dataset, MuHAVi dataset, IXMAS dataset

---

## 1. Introduction

Human action recognition has been of great interest in recent years due to its direct application and need in Surveillance, Ambient Intelligence, Ambient-Assisted Living (AAL) and Human-Computer Interaction systems. While it is still a recent field of research, huge advances have been made in classification of human actions (Poppe, 2010; Turaga et al., 2008; Weinland et al., 2011), recognition based on context and scene understanding (Kjellström, Sidenbladh; Bremond, 2007), as well as enhancement of traditional tracking and motion analysis systems with semantics about human activities (Moeslund et al., 2006; Hu et al., 2004). In this paper, a simple but yet very effective approach is presented in order to support accurate human action recognition at the level of basic human motion, like *walking, jumping, running, falling, etc.* Based on human silhouettes, a scale and location invariant feature is computed which shows to be a powerful discriminating signal, especially when considering its variation over time. At the training stage, the method learns the per class features that make up the most characteristic poses, the so called *key poses*. These can be acquired from single- or multi-view data, which makes the method suitable for scenarios with one or more cameras without any explicit constraints about the point of view (POV).

20 Using the ground truth data, the *sequences of key poses* corresponding to  
21 the labelled videos are obtained. These sequences are matched later with the  
22 current test sequence based on Dynamic Time Warping (DTW).

23 Our system has been designed so as to run at a frame rate close to real-  
24 time and to support online recognition. Since our target application is human  
25 monitoring at home for AAL services, these were both essential premises.  
26 Experimentation on three popular benchmarks (Weizmann from Blank et al.  
27 (2005), MuHAVi from Singh et al. (2010) and IXMAS from Moeslund et al.  
28 (2006)) shows that our approach outperforms state-of-the-art methods with  
29 similar conditions.

30 The contributions to the literature of this paper are two-fold. On the one  
31 hand, an efficient human action recognition method is presented which can  
32 be applied in a wide spectrum of application scenarios due to its performance  
33 in real-time and the absence of requirements as camera calibration or specific  
34 POVs. On the other hand, in this work human action recognition is carried  
35 out based on sequences of key poses. This achieves to filter noise and outliers  
36 from the training instances while at the same time it models the temporal  
37 evolution between key poses.

38 The remainder of this paper is organised as follows: section 2 summarises  
39 the most relevant and recent related works in human action recognition.  
40 In section 3 the chosen pose representation is analysed briefly. Our model  
41 learning approach is broken down into steps in section 4, and the final action  
42 recognition stage is presented in section 5. Section 6 gives a detailed analy-  
43 sis about the experimental results obtained and compares them with other  
44 state-of-the-art references. Finally, section 7 presents some conclusions and

45 discussion.

## 46 **2. Related Work**

47 When analysing human action recognition approaches based on vision  
48 techniques, classification can be made with respect to different semantic lev-  
49 els. Common criteria are: 1) the structural layout of the recognition method  
50 (Aggarwal and Ryoo, 2011); 2) the learning approach, for instance, exemplar-  
51 based vs. model-based, where we find generative models like Hidden Markov  
52 Models (HMM) and discriminative models like Conditional Random Fields  
53 (CRF) (Poppe, 2010); 3) the type of input features used for the classification  
54 (Poppe, 2010; Weinland et al., 2010).

55 Attending to the latter, *global* (also known as *dense* or *holistic*) represen-  
56 tations and *local* (also known as *sparse*) representations of the images can  
57 be obtained. The first require a region of interest (ROI) and therefore the  
58 human body needs to be detected in the image, usually with background  
59 subtraction and blob extraction techniques. While this additional step of  
60 pre-processing is a disadvantage, it is usually overcome by the significant  
61 reduction of both image size and inherent complexity of its content. Bobick  
62 and Davis (2001) used such a global representation in their Motion Histo-  
63 ry- and Energy-Images (MHI, MEI), which encode the temporal evolution  
64 of the movement of the image and its spatial location respectively over a  
65 sequence of frames. Weinland et al. (2006) extended the work of Bobick  
66 and Davis (2001) to a 3D Motion History Volume in order to combine im-  
67 ages from multiple cameras and to obtain a free-viewpoint representation.  
68 While Bobick and Davis (2001) use seven Hu Moments for description and

69 classification, Weinland et al. (2006) use Fourier analysis in cylindrical coor-  
70 dinates. Space-time volumes are constructed in Blank et al. (2005) by means  
71 of obtaining the solution to the Poisson equation for a sequence of binary  
72 silhouettes. Global space-time features (composed of the weighted moments  
73 of local space-time saliency and orientation features) are employed to achieve  
74 action recognition, detection and clustering. More recently, MHI templates  
75 have been clustered in a Self-Organising Map in order to represent image  
76 viewpoint and movement in a principal manifold (Martinez-Contreras et al.,  
77 2009). Each sequence of MHI is projected onto the map and the coordinates  
78 of activation are modelled with an HMM. Maximum Likelihood classifier is  
79 used for the final recognition.

80 There are also works which take advantage of image features that have  
81 not been originally designed for action recognition. Image gradients and op-  
82 tical flow have been widely and successfully used in tracking methods and  
83 their application to action recognition shows good results. In this sense, Tran  
84 and Sorokin (2008) designed a complex combination of shape and motion fea-  
85 tures. A 286-dimensional descriptor is obtained by encoding the binary shape  
86 of the silhouette, the vertical and horizontal optical flow and the context of  
87 15 surrounding frames reduced with PCA. Nearest Neighbour classification  
88 is done by discriminative metric learning and data subsampling. Fathi and  
89 Mori (2008) use mid-level motion features (spatio-temporal cuboids) made up  
90 of weighted combinations of thresholded low-level features based on optical  
91 flow. A variant of Adaboost is applied and one binary classifier is learned for  
92 every pair of classes in order to obtain a multi-class classifier, which achieves  
93 highly accurate results on popular action recognition datasets (Weizmann

94 from Blank et al. (2005) and KTH from Schuldt et al. (2004)). Main disad-  
95 vantages of such global representations are the lack of resistance to viewpoint  
96 changes and partial occlusions; under these circumstances global representa-  
97 tions suffer from high intra-class variance and are therefore difficult to learn  
98 accurately.

99 When using local representations, the image is regularly taken as it is  
100 and observed as a collection of patches or points. Commonly different types  
101 of salient points are obtained based on shape and gradient changes (like  
102 Harris and SUSAN corners, SIFT and SURF points; see Wu et al. (2010b);  
103 Juan and Gwun (2009) for more details). When considering the temporal  
104 evolution of the location or aspect of these points, space-time corners are  
105 applied. These encode 3D information of interest points “where the local  
106 neighbourhood has a significant variation in both the spatial and the tempo-  
107 ral domain” (Poppe, 2010). Great effort has been made to extend traditional  
108 salient point detectors to 3D: Laptev (2005) used the Harris corner as ba-  
109 sis, while Oikonomopoulos et al. (2005) extended the salient point detector  
110 from Kadir and Brady (2003), and Scovanner et al. (2007) created a 3D ver-  
111 sion of the popular SIFT points. A different approach is presented in İközler  
112 and Duygulu (2007), where the human body is represented with oriented  
113 rectangular patches; then a histogram is obtained with the  $15^\circ$  orientations  
114 resulting in 12 circular bins. Spatial information is encoded using a 3x3 grid  
115 and concatenating the histograms of each individual bin. Among different  
116 recognition methods, DTW showed the best results achieving perfect accu-  
117 racy with the Weizmann dataset. While local representations have achieved  
118 good recognition rates, great obstacles persist in attaining stable and con-

119 stant features in cluttered environments.

120 For greater detail about these methods and exhaustive reviews about the  
121 state of the art, we refer to the popular works Poppe (2010) and Moeslund  
122 et al. (2006), or more recent ones, like Aggarwal and Ryoo (2011); Chaaraoui  
123 et al. (2012).

### 124 3. Pose Representation

125 As introduced in section 1, our method relies on a global pose representa-  
126 tion based on the contour points of the silhouette. We assume that a binary  
127 silhouette is obtained previously by human silhouette extraction techniques,  
128 e.g. background subtraction. Using only the contour points and not the  
129 whole silhouette is motivated by getting rid of the redundancy that intro-  
130 duces the inside part of the human silhouette, leading therefore to a less  
131 expensive feature extraction. In addition, usage of contours avoids the need  
132 of morphological pre-processing steps and reduces the sensitivity to small  
133 viewpoint variations or lighting changes (Ángeles Mendoza and Pérez de la  
134 Blanca, 2007). Specifically, the contour-based feature from Dedeoğlu et al.  
135 (2006) has been chosen, which is described briefly in the following.

136 First, the contour points  $P = \{p_1, p_2, \dots, p_n\}$  of the silhouette need to be  
137 obtained. For this purpose, contour extraction is applied based on the border  
138 following algorithm from Suzuki and Be (1985).

139 Second, the centre of mass  $C_m = (x_c, y_c)$  of the silhouette's contour points  
140 is calculated with respect to the  $n$  number of points:

$$x_c = \frac{\sum_{i=1}^n x_i}{n}, y_c = \frac{\sum_{i=1}^n y_i}{n}. \quad (1)$$



141 Third, the distance signal  $DS = \{d_1, d_2, \dots, d_n\}$  is generated by deter-  
 142 mining the Euclidean distance between each contour point and the centre of  
 143 mass. Contour points should be considered always in the same order. For  
 144 instance, the set of points can start at the most left point with equal y-axis  
 145 value as the centre of mass, and follow a clockwise order.

$$d_i = \|C_m - p_i\|, \quad \forall i \in [1..n]. \quad (2)$$

146 Finally, scale-invariance is achieved by fixing the size of the distance sig-  
 147 nal, sub-sampling the feature size to a constant length  $L$ , and normalising  
 148 its values to unit sum.

$$\hat{DS}[i] = DS\left[i * \frac{n}{L}\right], \quad \forall i \in [1..L], \quad (3)$$

$$\bar{DS}[i] = \frac{\hat{DS}[i]}{\sum_{i=1}^L \hat{DS}[i]}, \quad \forall i \in [1..L]. \quad (4)$$

149 This type of global pose representation has a significant advantage over  
 150 similar features presented in section 2. While the spatial information is  
 151 preserved in greater detail than histogram- or grid-based representations,  
 152 the feature still has a low dimensionality and its processing presents a very  
 153 low computational cost (see section 6).

#### 154 4. Model Learning

155 Lately, several works (Baysal et al., 2010; Cheema et al., 2011; Eweiwi  
 156 et al., 2011; Thureau and Hlaváč, 2007) build upon key poses. Baysal et al.  
 157 (2010) define key poses as “a set of frames that uniquely distinguishes an

158 action from others”. Therefore, the goal of using key poses is to model an  
159 action by its most characteristic poses in time. This makes it possible to  
160 significantly reduce the problem scale in exemplar-based recognition meth-  
161 ods and, at the same time, to avoid redundant or superfluous learning. The  
162 underlying idea is that if the human brain is able to recognise what a person  
163 is doing based on a few individual images, why should not action recognition  
164 methods be able to sustain only on pose information. In this regard, Baysal  
165 et al. (2010); Cheema et al. (2011) use no temporal information at all, Thureau  
166 and Hlaváč (2007) model the short-term temporal relation between consec-  
167 utive key poses with n-grams (*trigrams* showed good results at acceptable  
168 computational cost), and Eweiwi et al. (2011) take into account the tempo-  
169 ral context of a small number of frames by means of obtaining temporal key  
170 poses based on MHI. While our approach is very similar to these works at  
171 the training stage when applied to a single view, our contribution considers  
172 long-term temporal relation between key poses and thus takes advantage of  
173 the known temporal evolution of key poses over a whole sequence.

174 A complete overview of the involved stages of the learning process can be  
175 seen in figure 1.

#### 176 4.1. Learning Key Poses

177 The first step of the learning process is to process all the frames of the  
178 video sequences in order to obtain their pose representation, as mentioned in  
179 section 3. Then, similar to Cheema et al. (2011); Baysal et al. (2010), the per  
180 class key poses are learned by means of *K-means* clustering with Euclidean  
181 distance. Hence, the extracted features of all available images of the same  
182 action class *samples* =  $\{s_1, s_2, \dots, s_n\}$  are grouped into  $K$  clusters; where

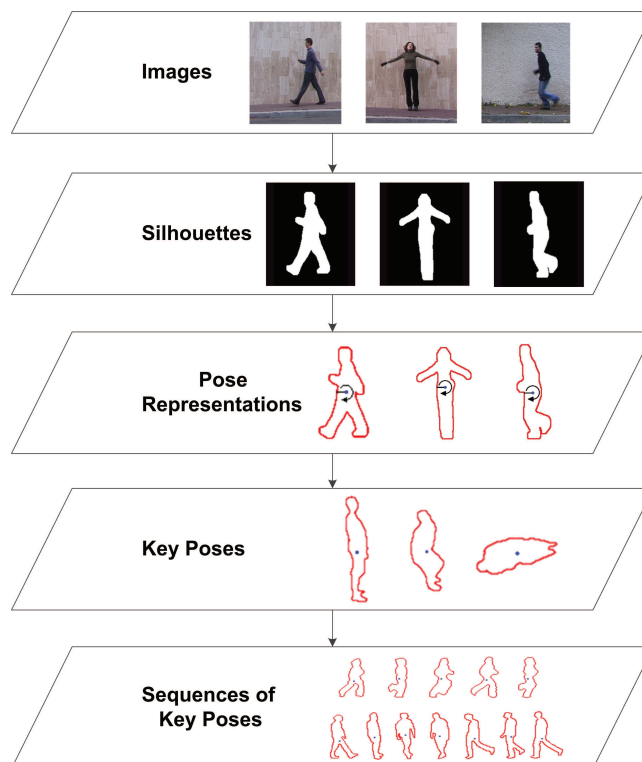


Figure 1: Overview of the learning process: first, a human silhouette extraction technique, like background subtraction, needs to be applied. Then the extracted human silhouettes are processed in order to obtain the contour-based feature. Finding the most characteristic poses among the training data returns the key poses. The sequences of key poses model the temporal evolution between key poses with respect to the original training sequences.

183 each cluster centre of  $centres = \{c_1, c_2, \dots, c_K\}$  represents a key pose  $kp$  as it  
 184 is a characteristic pose among the training data. The process of clustering is  
 185 repeated  $\lambda$  times, so as to avoid local minimum, and the best result is taken  
 186 (the usage of more advanced clustering algorithms is being considered for  
 187 future works). Given that the clustering process returns the corresponding  
 188 label of each sample,  $labels = \{l_1, l_2, \dots, l_n\}$  in which  $l_i$  stands for the index of  
 189 the cluster assigned to  $s_i$ , clustering results are evaluated with the following  
 190 compactness metric  $C$ :

$$C = \sum_{i=1}^n |s_i - c_{l_i}|, \quad (5)$$

191 where the instance with the lowest value is taken as the final result.

192 This key pose learning process is repeated individually for the training  
 193 samples of each action class. This way, a set of  $K$  key poses is obtained for  
 194 each action class.

#### 195 4.2. Learning Sequences of Key Poses

196 As stated beforehand, our goal is to learn the long-term temporal evo-  
 197 lution of key poses. Consequently, our interest resides on the successive  
 198 key poses that are involved in an action performance. As the training data  
 199 is made up of sequences of labelled action performances, the correspond-  
 200 ing sequences of key poses can be modelled. For the pose representation of  
 201 each frame of a sequence, i.e.  $S_{poses} = \{pose_1, pose_2, \dots, pose_n\}$ , the *nearest*  
 202 *neighbour* key pose is found. The successive *nearest neighbour* key poses  
 203 constitute the simplified sequence of known characteristic poses and their  
 204 evolution:  $S = \{kp_1, kp_2, \dots, kp_n\}$ . This way, a set of sequences of key poses

205 is obtained for each action class. This decisive step significantly improves  
206 exemplar-based action recognition by shifting the training data to a com-  
207 mon and known domain (the set of characteristic key poses), and therefore  
208 filtering out single examples with noise or partial occlusions.

### 209 4.3. Learning from Multiple Views

210 Nowadays, most application scenarios do have more than one camera  
211 available. Multiple views of the same environment help to avoid occlusions  
212 due to obstacles (like furniture or having several persons in the field of view)  
213 and make it possible to have multiple POV of the same event at our disposal.  
214 However, the task of dealing with several video streams, modelling 3D repre-  
215 sentations and targeting action recognition applications still has to overcome  
216 great difficulties, as dealing with richer data leads to high computational cost  
217 and burdensome systems (Moeslund et al., 2006; Holte et al., 2011).

218 Since the presented method shows successful results in single-view action  
219 recognition, one wonders if the approach is able to accurately model multi-  
220 view data. Among the different available approaches of combining multi-view  
221 data (Holte et al., 2011; Wu et al., 2010a) a *feature fusion* approach has been  
222 chosen, so as to test if the model based on sequences of key poses is able to  
223 learn from multiple views. In this sense, multi-view data is combined at the  
224 feature level and no changes are performed at the modelling or recognition  
225 levels.

226 Assuming that  $v$  video streams of the same scenario are available, first  
227 each frame is individually processed to its pose representation. Then the  
228 multi-view pose representation  $\bar{D}S_{mv}$  is obtained by frame-by-frame con-

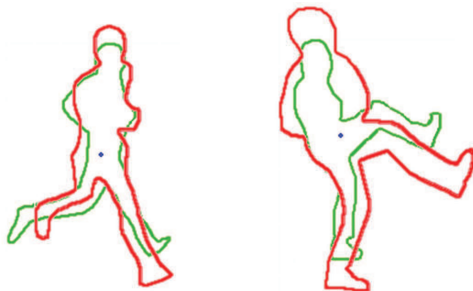


Figure 2: Multi-view key poses: *RunLeftToRight* (left) and *KickRight* (right) from MuHAVi.

229 catenation of single-view pose representations  $\bar{D}S_{sv}$ :

$$\bar{D}S_{mv} = \bar{D}S_{sv_1} \circ \bar{D}S_{sv_2} \circ \dots \circ \bar{D}S_{sv_v}. \quad (6)$$

230 This step is identically performed with train and test instances, using  
 231 multi-view pose representations at the succeeding stages. As a result, when  
 232 feeding the model with multi-view pose representations, sequences of multi-  
 233 view key poses (see figure 2) are inherently obtained.

## 234 5. Action Recognition

235 At the recognition stage, a final class label output needs to be given. To  
 236 that end, two steps have to be taken: 1) in the same way as with our training  
 237 sequences, silhouette contour points are processed and their corresponding  
 238 pose representations are obtained; 2) for each test sequence, the pose repre-  
 239 sentation of each frame is used to find the *nearest neighbour* key pose and  
 240 build the analogous sequence of *nearest neighbour* key poses. This shift to  
 241 our known data domain acts as filtering and simplification process, and in-  
 242 troduces the needed stability when dealing with test data with meaningful

243 differences to the training data, like action performances of different actors  
 244 (see section 6).

245 Due to the temporal intra-class variance, a suitable distance metric is  
 246 needed in order to compare the sequences of key poses. Different actors can  
 247 perform the same actions on very different ways and they can do so faster or  
 248 slower than others. While some motions are indispensable when performing  
 249 an action, like moving one leg and then the other while walking, these can still  
 250 appear with a considerable time shift, especially when dealing with elderly  
 251 people. Dynamic Time Warping is particularly suitable when dealing with  
 252 the comparison of sequences that can present inconsistent time scales, but  
 253 without changing the temporal order. It is able to align two time series of  
 254 different lengths even if there are accelerations or decelerations.

255 Given two sequences of key poses  $S_{train} = \{kp_1, kp_2, \dots, kp_n\}$  and  $S_{test} =$   
 256  $\{kp'_1, kp'_2, \dots, kp'_m\}$  we compute the DTW distance  $S_{train} - S_{test}$  as:

$$S_{train} - S_{test} = dtw(n, m), \quad (7)$$

$$dtw(i, j) = \min \left\{ \begin{array}{l} dtw(i-1, j), \\ dtw(i, j-1), \\ dtw(i-1, j-1) \end{array} \right\} + d(kp_i, kp'_j), \quad (8)$$

257 where  $d(kp_i, kp'_j)$  is the Euclidean distance used for feature comparison  
 258 between two key poses.

259 This way, using DTW, the nearest neighbour sequence of key poses is  
 260 found and its label supplies the final result.

## 261 6. Experimentation

262 In order to test the accuracy and stability of the presented approach,  
263 three human action recognition datasets have been used as benchmarks. In  
264 the case of the Weizmann dataset, a *leave-one-sequence-out* cross validation  
265 procedure has been applied. This way, the system is trained with all but one  
266 video sequence, which is the one that evaluates the accuracy score. Iterating  
267 over all the sequences, the average success rate is used as final result. In the  
268 case of the MuHAVi dataset, its authors introduced an evaluation scheme  
269 based on view- and actor-invariance tests which we repeat so as to compare  
270 our results. And in the IXMAS dataset we used the usual *leave-one-actor-out*  
271 cross validation. Finally, a temporal evaluation is made in order to confirm  
272 the suitability for real-time applications. A comparison of the presented  
273 results with similar state-of-the-art approaches is given in section 6.5.

274 The three constant parameters of the presented method have been chosen  
275 based on empirical testing. The number of clustering attempts  $\lambda = 3$  for  
276 all results shown, while the length of the distance signal feature  $L$  and the  
277 number of key poses per action class  $K$  are detailed for each test.

### 278 6.1. Weizmann Dataset

279 The Weizmann dataset presented in Blank et al. (2005) is a single-view  
280 (static front-side camera) outdoor dataset. It provides 180x144 px resolution  
281 images of 10 different actions performed by 9 actors. It has a relatively  
282 simple background, provides automatically extracted silhouettes (we use the  
283 version without post-alignment), and has become a reference in human action  
284 recognition. Actions include *bending (bend)*, *jumping jack (jack)*, *jumping*



	bend	jack	jump	pjump	run	side	walk	wave1	wave2
bend	9/9								
jack		9/9							
jump			9/9						
pjump				9/9					
run					7/10		3/10		
side			1/9			8/9			
walk							10/10		
wave1		1/9						8/9	
wave2		1/9							8/9

Figure 3: Confusion matrix of the Weizmann dataset without the *skip* action. *Leave-one-sequence-out* cross validation with 83 sequences.

285 *forward (jump), jumping in place (pjump), running (run), galloping sideways*  
286 *(side), skipping (skip), walking (walk), waving one hand (wave1) and waving*  
287 *two hands (wave2).* It is worth mentioning that several works exclude the  
288 *skip* action, as it commonly shows higher error rates and also weakens the  
289 recognition of other actions.

290 Figure 3 shows the result of the cross validation test without the *skip*  
291 action. At an average success rate of 92.77% (achieved with  $L = 120$  and  
292  $K = 96$ ), it can be seen that the confusions made are coherent. As seen in  
293 the works from Saghafi and Rajan (2012); Shao and Chen (2010), *walk* and  
294 *run* present a high inter-class similarity, and therefore the difference between  
295 their key poses is minimal. In *jack* hands are risen, similarly to *wave1* and  
296 *wave2*.

297 Taking a closer look to the misclassifications of sequences from the *run*  
298 action class, it can be seen that the running or walking speed of the ac-  
299 tors varied significantly. In addition, some of the actors do not move their  
300 arms along when running, which increases even more the similarity between  
301 running and walking. We have analysed a specific misclassification of a *run*

Table 1: Ten closest key pose sequences for a specific misclassification of a *run* sequence.

Index	Action class	DTW distance
1	walk	3,264716
2	run	3,795877
3	walk	4,116315
4	side	4,722770
5	run	4,869563
6	run	5,224457
7	run	5,319681
8	run	5,458966
9	run	6,019087
10	run	6,206304

302 sequence (see table 1). The ten closest sequences include seven sequences  
 303 of the right class, which means that, for instance, a K-Nearest Neighbour  
 304 (KNN) approach could have worked better in this case. The sequence num-  
 305 ber 2 is the closest sequence that would have produced a successful match.  
 306 A 100% of its key poses proceed from the training instances of the *run* class.  
 307 Surprisingly, only  $\sim 14\%$  of the frames of the tested sequence have matched  
 308 with a key pose from this class, which explains why this sequence has been  
 309 misclassified.

310 When including the *skip* action, the success rate decreases to 90.32%  
 311 (achieved with  $L = 200$  and  $K = 96$ ). Interestingly, this action is recognised  
 312 perfectly, but the stability of the other actions is still affected because of the  
 313 rise of inter-class similarity which occurs when adding this action class. It  
 314 has been observed that the *skip* key poses get hit very frequently in several  
 315 action classes as *jump*, *pjump*, *run*, *side* and *walk*. Similar conclusions have

316 been obtained in Saghafi and Rajan (2012); Shao and Chen (2010).

## 317 6.2. *MuHAVi Dataset*

318 The MuHAVi dataset (Singh et al., 2010) is a more recent and com-  
319 plex benchmark with multi-view images. It provides 720x576 px resolu-  
320 tion images on a complex background with street light illumination. Its  
321 full version includes 17 different actions performed by 7 actors and has been  
322 recorded indoors with 8 CCTV cameras, each one at 45° to its neighbours.  
323 A manually annotated subset (*MuHAVi-MAS*) provides silhouettes for 2  
324 of these views (front-side and 45°) and 2 actors, labelling 14 (MuHAVi-  
325 14: *CollapseLeft*, *CollapseRight*, *GuardToKick*, *GuardToPunch*, *KickRight*,  
326 *PunchRight*, *RunLeftToRight*, *RunRightToLeft*, *StandupLeft*, *StandupRight*,  
327 *TurnBackLeft*, *TurnBackRight*, *WalkLeftToRight* and *WalkRightToLeft*) or 8  
328 (MuHAVi-8: *Collapse*, *Guard*, *KickRight*, *PunchRight*, *Run*, *Standup*, *Turn-*  
329 *Back* and *Walk*) actions in its merged version.

### 330 6.2.1. *Leave-one-sequence-out Cross Validation*

331 As this dataset includes multi-view data, our method uses the proposed  
332 multi-view pose representations and learns sequences of multi-view key poses.  
333 Since two camera views are available, sequences are considered as pairs, each  
334 of which contains the images of the same action performance from a differ-  
335 ent view. Therefore, the 136 available sequences are taken as 68 different  
336 sequences when performing the *leave-one-sequence-out* cross validation test.

337 In figure 4, the confusion matrix for MuHAVi-14 shows very promising  
338 results with an average success rate of 91.18% (achieved with  $L = 340$  and  
339  $K = 90$ ), misclassifying only 6 sequences.

	CollapseLeft	CollapseRight	GuardToKick	GuardToPunch	KickRight	PunchRight	RunLeftToRight	RunRightToLeft	StandupLeft	StandupRight	TurnBackLeft	TurnBackRight	WalkLeftToRight	WalkRightToLeft
CollapseLeft	4/4													
CollapseRight		4/4												
GuardToKick			6/8	2/8										
GuardToPunch				8/8										
KickRight					8/8									
PunchRight						8/8								
RunLeftToRight							3/4	1/4						
RunRightToLeft								4/4						
StandupLeft									1/2	1/2				
StandupRight										4/4				
TurnBackLeft			1/2								1/2			
TurnBackRight			1/4									3/4		
WalkLeftToRight													4/4	
WalkRightToLeft														4/4

	Collapse	Guard	KickRight	PunchRight	Run	Standup	TurnBack	Walk
Collapse	8/8							
Guard		16/16						
KickRight			8/8					
PunchRight				8/8				
Run					8/8			
Standup						6/6		
TurnBack		2/6					4/6	
Walk								8/8

Figure 4: Confusion matrices of the MuHAVi dataset: MuHAVi-14 (top) and MuHAVi-8 (bottom). *Leave-one-sequence-out* cross validation with 68 multi-view sequences.

340 In MuHAVi-8 only 2 sequences are misclassified and a success rate of  
341 97.06% ( $L = 250$  and  $K = 90$ ) is achieved. In both tests it can be seen that  
342 *TurnBack* shows greater difficulty than other actions.

### 343 6.2.2. Identical Actors, Novel Camera

344 In this view-invariance test, all available sequences of one POV are used  
345 at training, whereas at testing, the same sequences but from the second  
346 POV are used. Hence, no multi-view learning can be applied. This test is  
347 executed twice, interchanging the training and testing groups, and the results  
348 are averaged.

349 Since view-invariance has not been explicitly considered, no exceptional  
350 robustness is expected in this sense. The test returns a result of 38.97%  
351 ( $L = 220$  and  $K = 70$ ) on MuHAVi-14 and 63.24% ( $L = 370$  and  $K = 50$ )  
352 on MuHAVi-8.

### 353 6.2.3. Identical Cameras, Novel Actor

354 Similarly to the last test, all sequences of one actor are used at training,  
355 while the sequences of a different actor, unknown to the learning model, are  
356 used at testing (and vice-versa). As more than one view of the same action  
357 performance is available, multi-view learning is applied and 34 sequences  
358 with images of two views are used at training and another 34 at testing.

359 In contrast to the last test and as mentioned before, the presented method  
360 is designed to be robust to test data with meaningful differences to the train  
361 data (due to dissimilarities among actors or noise). For this reason, data  
362 is first shifted to the known domain of key poses and then matched to the  
363 corresponding train sequence.

364 Actor-invariance tests present an increased difficulty due to the singular-  
365 ity of multiple actor-dependant conditions. In this sense, parameters as size,  
366 body build, clothes, etc. are given by the actor, as well as the particular  
367 way in which each person performs an action. This can be seen, for instance,  
368 in gait analysis, where the involved dynamics even allow to perform person  
369 identification (Wang et al., 2010).

370 The *Novel Actor* test returns a success rate of 82.35% ( $L = 450$  and  
371  $K = 110$ ) on MuHAVi-14 and 88.24% ( $L = 250$  and  $K = 110$ ) on MuHAVi-  
372 8. To the best of our knowledge, these are the highest results achieved so  
373 far.

### 374 6.3. IXMAS Dataset

375 With the purpose of extending the experimentation of our method to  
376 a more difficult dataset with more camera views, we have chosen the IX-  
377 MAS dataset which is popular among human action recognition methods  
378 that are specifically designed for multi-view recognition. The INRIA Xmas  
379 Motion Acquisition Sequences (IXMAS) dataset (Weinland et al., 2006) in-  
380 cludes multi-view data and is especially aimed at view-invariance testing. It  
381 provides 390x291 px resolution images from five different angles including  
382 four sides and one top-view camera. A set of 12 actors have been recorded  
383 performing 14 different actions (*check watch, cross arms, scratch head, sit*  
384 *down, get up, turn around, walk, wave, punch, kick, point, pick up, throw*  
385 *over head* and *throw from bottom up*) 3 times each, resulting in a dataset  
386 with over 2 000 sequences. This benchmark presents an increased difficulty  
387 because subjects were asked to freely choose their position and orientation.  
388 Therefore, each camera has captured different viewing angles, which makes

	check watch	cross arms	scratch head	sit down	get up	turn around	walk	wave	punch	kick	pick up
check watch	28/36	4/36							4/36		
cross arms		32/36	1/36						3/36		
scratch head	1/36	8/36	22/36						2/36	2/36	1/36
sit down				35/36							1/36
get up					36/36						
turn around		1/36				35/36					
walk						6/36	30/36				
wave	3/36		6/36					26/36	1/36		
punch	3/36	1/36	1/36					1/36	30/36		
kick		1/36							4/36	30/36	1/36
pick up											36/36

Figure 5: Confusion matrix of the IXMAS dataset. *Leave-one-actor-out* cross validation with 11 actors and 396 multi-view sequences.

389 methods which rely on fixed camera views (front, side, etc.) unsuitable.

390 Figure 5 shows the confusion matrix that has been obtained for this chal-  
391 lenging dataset. As common in the state-of-the-art, we used a *leave-one-*  
392 *actor-out* cross validation test in which actor-invariance is tested by training  
393 with the instances from all but one actor and testing the sequences from  
394 the unknown one. This is repeated for all available actors and the average  
395 accuracy score is obtained. Following the test setup given by the publishers  
396 of the dataset, we excluded the *point* and *throw* actions. The test returns  
397 an average result of 85.86% ( $L = 400$  and  $K = 20$ ). As it can be seen in  
398 the confusion matrix, the actions that are performed with arms and hands  
399 present several misclassifications due to their similarity. *Walk* is matched  
400 with *turn around* because the proposed method does only rely on silhouette  
401 shape without explicitly learning action’s kinematics. Turning around is es-  
402 sentially walking with a specific direction and this is not differentiated by  
403 our system.

#### 404 6.4. Temporal Evaluation

405 When designing a human action recognition method intended to perform  
406 online, the temporal constraint is crucial. Even more when considering that  
407 this unit would be only one part of a complex distributed vision system which  
408 performs movement detection, tracking, background segmentation, person  
409 identification, privacy filtering, etc., and moreover needs to be executed on  
410 an embedded hardware device. For this reason, a human action recognition  
411 module needs to perform as fast as possible, and simple yet effective ap-  
412 proaches are preferred over perfect yet unaffordable ones. Our evaluation  
413 system consists of a standard PC with an Intel Core 2 Duo CPU at 3 GHz,  
414 running Windows 7 64-bit and an implementation using the .NET Frame-  
415 work and the widely used Computer Vision library OpenCV (Bradski, 2000).  
416 Time evaluation has been performed using the hardware counter *QueryPer-*  
417 *formanceCounter* with a precision of  $\mu s$ .

418 Executing the learning process for the 93 sequences of the Weizmann  
419 dataset, which contain 5687 frames of 180x144 px, takes 81.1s. That is an  
420 average of 0.87s per sequence at 70.12FPS. But more important is the speed  
421 of the testing process which takes 45.72s, achieving an average speed of 0.49s  
422 per sequence at 124.38FPS.

423 In MuHAVi-14, the training of 136 sequences made up of 7941 frames  
424 of 720x576 px takes 204.44s, i.e. an average speed of 1.5s per sequence  
425 at 38.84FPS. The testing process for this data takes 109.9s, achieving an  
426 average speed of 0.81s per sequence at 72.25FPS. As MuHAVi-8 has fewer  
427 action classes, the learning process speeds up to 53.76FPS and the testing  
428 process to 81.31FPS.



Table 2: Comparison with similar state-of-the-art approaches on the Weizmann dataset.

Approach	Input	Actions	Evaluation	Rate	FPS
İkizler and Duygulu (2007)	Silhouettes	9	LOSO	100%	N/A
Tran and Sorokin (2008)	Silhouettes	10	LOSO	100%	N/A
Eweiwi et al. (2011)	Aligned sil.	10	LOSO	100%	N/A
Hernández et al. (2011)	Images	10	LOAO	90.3%	98
Cheema et al. (2011)	Silhouettes	9	LOSO	91.6%	56
Our method	Silhouettes	9	LOSO	92.8%	124

429 In the case of the IXMAS dataset these rates change to  $155.52FPS$  for  
 430 the training process and  $26.48FPS$  for the testing process.

431 These tests were performed including all processing stages from the com-  
 432 puting of the contour points to the actual recognition, and using the sil-  
 433 houette images as basis. The obtained performances correspond to the best  
 434 test configurations shown in previous sections, without applying any further  
 435 optimisation.

### 436 6.5. Comparison of Results

437 The comparison of different human action recognition approaches can  
 438 be difficult and misleading because of diverse recognition goals (some only  
 439 seek an action class label, and others need a reconstructed 3D environment),  
 440 different kinds of input data (images, video streams, silhouettes, outputs of  
 441 tracking systems, etc.) and even incompatible evaluation methods.

442 Table 2 shows a comparison of our result on the Weizmann dataset with  
 443 other similar approaches. The success rates are obtained either with *leave-*  
 444 *one-actor-out* (LOAO) or *leave-one-sequence-out* (LOSO) cross validations.  
 445 Several works achieve perfect recognition on this dataset, but most of them

Table 3: Comparison with similar state-of-the-art approaches on the MuHAVi dataset. All use silhouettes as input data and LOSO as evaluation method.

<b>Approach</b>	<b>MuHAVi-14</b>		<b>MuHAVi-8</b>	
	<b>Rate</b>	<b>FPS</b>	<b>Rate</b>	<b>FPS</b>
Singh et al. (2010) (baseline)	82.4%	N/A	97.8%	N/A
Martinez-Contreras et al. (2009)	-	-	98.4%	N/A
Eweiwi et al. (2011)	91.9%	N/A	98.5%	N/A
Cheema et al. (2011)	86.0%	56	95.6%	56
Our method	91.2%	72	97.1%	81

Table 4: Comparison of results of the MuHAVi Novel Actor test.

<b>Approach</b>	<b>MuHAVi-14</b>	<b>MuHAVi-8</b>
Singh et al. (2010)	61.8%	76.4%
Cheema et al. (2011)	73.5%	83.1%
Eweiwi et al. (2011)	77.9%	85.3%
Our method	82.4%	88.2%

446 do not present any temporal evaluation and their suitability for real-time  
 447 applications is arguable. It can be seen that, when comparing with methods  
 448 that present temporal data, our performance improves state-of-the-art rates  
 449 both in recognition accuracy and speed.

450 Table 3 presents similar comparisons for the MuHAVi dataset. Again  
 451 the present method achieves state-of-the-art success rates and outperforms  
 452 similar methods with real-time suitability in recognition accuracy, as well as  
 453 in recognition speed.

454 We also want to point out the robustness of our method with respect  
 455 to the *Novel Actor* test. Dissimilarities among action performances from

Table 5: Comparison with other multi-view human action recognition approaches of the state-of-the-art. The rates obtained in the *leave-one-actor-out* cross validation performed on the IXMAS dataset are shown (except for Cherla et al. (2008) where the type of test is not stated).

<b>Approach</b>	<b>Input</b>	<b>Actions</b>	<b>Actors</b>	<b>Views</b>	<b>Rate</b>	<b>FPS</b>
Wu et al. (2011)	Images	12	12	4	89.4%	N/A
Weinland et al. (2006)	Silhouettes	11	10	5	93.3%	N/A
Holte et al. (2012)	Images	13	12	5	100%	N/A
Cherla et al. (2008)	Silhouettes	13	N/A	4	80.1%	20
Weinland et al. (2010)	Images	11	10	5	83.5%	~500
Our method	Silhouettes	11	12	5	85.9%	26

456 different actors lie in speed, shape and motion. As shown in table 4, our  
 457 approach clearly outperforms latest results on both versions of the MuHAVi  
 458 dataset. As seen in the results from Singh et al. (2010) and Cheema et al.  
 459 (2011), this test presents a higher difficulty and the improvements achieved  
 460 by our proposal constitute a significant benefit.

461 Last but not least, we compared the results obtained on the IXMAS  
 462 dataset which presented a much higher degree of difficulty due to its increased  
 463 number of actions, actors and views, as well as the different orientations that  
 464 the subjects chose with respect to the cameras. Table 5 shows a comparison  
 465 with other multi-view human action recognition approaches. The number  
 466 of action classes, actors and views have been detailed because these vary  
 467 among the approaches. Wu et al. (2011) obtained their highest rate excluding  
 468 camera 4, whereas Cherla et al. (2008) excluded the top-view camera and  
 469 reorganised the 4 side views into 6 viewing angles in order to achieve view  
 470 consistency. Recently, Holte et al. (2012) achieved perfect recognition on

471 this dataset relying on 4D spatio-temporal interest points. Nonetheless, the  
472 published recognition rates decrease when searching for methods which prove  
473 to be suitable for real-time applications. Once again, our method shows to  
474 be superior when regarding both action recognition accuracy and speed.

475 It can be seen that the improvements achieved for the MuHAVi dataset  
476 are more significant, and this is directly related to the quality of the input  
477 data. The silhouettes from the Weizmann and IXMAS datasets have been au-  
478 tomatically extracted through background subtraction techniques. For this  
479 reason, the results present noise and incompleteness. Although, real-time  
480 silhouette extraction of an acceptable quality can be performed (Horprasert  
481 et al., 1999; Kim et al., 2005), silhouettes of a substantial higher quality  
482 can be obtained by recent advances in depth sensors which are able to ap-  
483 ply markerless human pose recognition in real-time (Shotton et al., 2011).  
484 Furthermore, as the employed feature relies on the raw contour data and  
485 therefore presents sensitivity to these type of errors, image filters as border  
486 smoothing could be applied; or a more robust feature proposal could be used.

## 487 **7. Conclusion and Discussion**

488 In this paper, we have presented a human action recognition approach  
489 based on sequences of key poses. The human silhouette obtained, for in-  
490 stance, with background subtraction is used as initial input. The silhouette's  
491 contour leads to the used pose representation, by means of a distance sig-  
492 nal feature which, in conjunction with the model learning approach and the  
493 action classification, shows to be a highly efficient technique. Accurate recog-  
494 nition results are obtained without compromising the method's suitability for

495 real-time applications.

496 In contrast to exemplar-based methods, choosing a key pose-based ap-  
497 proach leads to a simplified classification process in which the number of ref-  
498 erence patterns is drastically reduced and noisy examples are filtered. The  
499 sequences of key poses allow us to model the long-term temporal evolution  
500 involved in action performances. Since the key poses themselves are non-  
501 temporal, introducing the temporal relationship between them at a supe-  
502 rior level allows a higher semantic richness and improves classification with  
503 respect to strictly non-temporal key pose-based methods. Finally, an ap-  
504 propriate and efficient sequence matching algorithm, like DTW, enables to  
505 successfully classify sequences with inconsistent time scales. As section 6  
506 shows, the presented method returns highly promising results on publicly  
507 available datasets, deals with both single- and multi-view scenarios success-  
508 fully, and is especially robust to different ways in which actions are performed  
509 by different actors.

510 However, when considering sequences of key poses, we assume that the  
511 temporal order is always the same, limitation that could be overcome with  
512 the use of probabilistic graphical models like HMM. Moreover, as our method  
513 does not take into account location or optical flow, the system would have  
514 difficulty in distinguishing, for instance, walking forwards from walking back-  
515 wards, because the involved poses and their relation are nearly identical.  
516 Other future lines include evaluating our method using images with occlu-  
517 sions and recognising a *null* or *unknown* action class which defines the normal  
518 human behaviour. The latter could be classified based on the distances to  
519 the learned action classes. If none of them is a good match, the *unknown*

520 action class can be hit. Finally, view-invariance is not taken into account  
521 and different subject orientations need to be learned explicitly.

### 522 *Acknowledgements.*

523 This work has been partially supported by the Spanish Ministry of Sci-  
524 ence and Innovation under project “Sistema de visión para la monitorización  
525 de la actividad de la vida diaria en el hogar” (TIN2010-20510-C04-02) and  
526 by the European Commission under project “caring4U - A study on people  
527 activity in private spaces: towards a multisensor network that meets pri-  
528 vacy requirements” (PIEF-GA-2010-274649). Alexandros Andre Chaaraoui  
529 acknowledges financial support by the Conselleria d’Educació, Formació i  
530 Ocupació of the Generalitat Valenciana (fellowship ACIF/2011/160). The  
531 funders had no role in study design, data collection and analysis, decision to  
532 publish, or preparation of the manuscript.

### 533 **References**

534 Aggarwal, J., Ryoo, M., 2011. Human activity analysis: A review. ACM  
535 Comput. Surv. 43, 16:1–16:43.

536 Baysal, S., Kurt, M., Duygulu, P., 2010. Recognizing human actions us-  
537 ing key poses, in: Pattern Recognition (ICPR), 2010 20th International  
538 Conference on, pp. 1727 –1730.

539 Blank, M., Gorelick, L., Shechtman, E., Irani, M., Basri, R., 2005. Actions  
540 as space-time shapes, in: Computer Vision, 2005. ICCV 2005. Tenth IEEE  
541 International Conference on, pp. 1395 –1402 Vol. 2.

- 542 Bobick, A., Davis, J., 2001. The recognition of human movement using  
543 temporal templates. *Pattern Analysis and Machine Intelligence*, IEEE  
544 *Transactions on* 23, 257 –267.
- 545 Bradski, G., 2000. The OpenCV Library. *Dr. Dobb's Journal of Software*  
546 *Tools* .
- 547 Bremond, F., 2007. Scene Understanding: perception, multi- sensor fusion,  
548 spatio-temporal reasoning and activity recognition. Ph.D. thesis. Univer-  
549 sité de Nice-Sophia Antipolis.
- 550 Chaaaraoui, A.A., Climent-Pérez, P., Flórez-Revuelta, F., 2012. A review on  
551 vision techniques applied to human behaviour analysis for ambient-assisted  
552 living. *Expert Systems with Applications* 39, 10873 – 10888.
- 553 Cheema, S., Eweiwi, A., Thureau, C., Bauckhage, C., 2011. Action recogni-  
554 tion by learning discriminative key poses, in: *Computer Vision Workshops*  
555 (*ICCV Workshops*), 2011 IEEE International Conference on, pp. 1302 –  
556 1309.
- 557 Cherla, S., Kulkarni, K., Kale, A., Ramasubramanian, V., 2008. Towards  
558 fast, view-invariant human action recognition, in: *Computer Vision and*  
559 *Pattern Recognition Workshops*, 2008. CVPRW '08. IEEE Computer So-  
560 ciety Conference on, pp. 1 –8.
- 561 Dedeoğlu, Y., Töreyn, B., Güdükbay, U., Çetin, A., 2006. Silhouette-based  
562 method for object classification and human action recognition in video,  
563 in: Huang, T., Sebe, N., Lew, M., Pavlovic, V., Kölsch, M., Galata, A.,  
564 Kisacanın, B. (Eds.), *Computer Vision in Human-Computer Interaction*.

- 565 Springer Berlin / Heidelberg. volume 3979 of *Lecture Notes in Computer*  
566 *Science*, pp. 64–77.
- 567 Eweiwi, A., Cheema, S., Thureau, C., Bauckhage, C., 2011. Temporal key  
568 poses for human action recognition, in: *Computer Vision Workshops*  
569 *(ICCV Workshops)*, 2011 IEEE International Conference on, pp. 1310 –  
570 1317.
- 571 Fathi, A., Mori, G., 2008. Action recognition by learning mid-level motion  
572 features, in: *Computer Vision and Pattern Recognition, 2008. CVPR 2008.*  
573 *IEEE Conference on*, pp. 1 –8.
- 574 Hernández, J., Montemayor, A., José Pantrigo, J., Sánchez, A., 2011. Human  
575 action recognition based on tracking features, in: Ferrández, J., Álvarez  
576 Sánchez, J., de la Paz, F., Toledo, F. (Eds.), *Foundations on Natural*  
577 *and Artificial Computation*. Springer Berlin / Heidelberg. volume 6686 of  
578 *Lecture Notes in Computer Science*, pp. 471–480.
- 579 Holte, M., Chakraborty, B., Gonzalez, J., Moeslund, T., 2012. A local 3-d  
580 motion descriptor for multi-view human action recognition from 4-d spatio-  
581 temporal interest points. *Selected Topics in Signal Processing, IEEE Jour-*  
582 *nal of* 6, 553–565.
- 583 Holte, M.B., Tran, C., Trivedi, M.M., Moeslund, T.B., 2011. Human ac-  
584 tion recognition using multiple views: a comparative perspective on recent  
585 developments, in: *Proceedings of the 2011 joint ACM workshop on Hu-*  
586 *man gesture and behavior understanding*, ACM, New York, NY, USA. pp.  
587 47–52.



- 588 Horprasert, T., Harwood, D., Davis, L., 1999. A statistical approach for  
589 real-time robust background subtraction and shadow detection, in: IEEE  
590 ICCV, pp. 256–261.
- 591 Hu, W., Tan, T., Wang, L., Maybank, S., 2004. A survey on visual surveil-  
592 lance of object motion and behaviors. *Systems, Man, and Cybernetics,*  
593 *Part C: Applications and Reviews*, IEEE Transactions on 34, 334 –352.
- 594 İkizler, N., Duygulu, P., 2007. Human action recognition using distribution  
595 of oriented rectangular patches, in: Elgammal, A., Rosenhahn, B., Klette,  
596 R. (Eds.), *Human Motion Understanding, Modeling, Capture and An-*  
597 *imation*. Springer Berlin / Heidelberg. volume 4814 of *Lecture Notes in*  
598 *Computer Science*, pp. 271–284.
- 599 Juan, L., Gwun, O., 2009. A Comparison of SIFT , PCA-SIFT and SURF.  
600 *International Journal of Image Processing (IJIP)* 3, 143 – 152.
- 601 Kadir, T., Brady, M., 2003. Scale saliency: a novel approach to salient feature  
602 and scale selection, in: *Visual Information Engineering, 2003. VIE 2003.*  
603 *International Conference on*, pp. 25 – 28.
- 604 Kim, K., Chalidabhongse, T.H., Harwood, D., Davis, L., 2005. Real-time  
605 foreground-background segmentation using codebook model. *Real-Time*  
606 *Imaging* 11, 172 – 185. Special Issue on Video Object Processing.
- 607 Kjellström (Sidenbladh), H., 2011. Contextual action recognition, in: Moes-  
608 lund, T.B., Hilton, A., Krüger, V., Sigal, L. (Eds.), *Visual Analysis of*  
609 *Humans*. Springer London, pp. 355–376.

- 610 Laptev, I., 2005. On space-time interest points. *International Journal of*  
611 *Computer Vision* 64, 107–123.
- 612 Martínez-Contreras, F., Orrite-Urunuela, C., Herrero-Jaraba, E., Ragheb,  
613 H., Velastin, S., 2009. Recognizing human actions using silhouette-based  
614 hmm, in: *Advanced Video and Signal Based Surveillance, 2009. AVSS '09.*  
615 *Sixth IEEE International Conference on*, pp. 43 –48.
- 616 Ángeles Mendoza, M., Pérez de la Blanca, N., 2007. Hmm-based action  
617 recognition using contour histograms, in: Martí, J., Benedí, J., Mendonça,  
618 A., Serrat, J. (Eds.), *Pattern Recognition and Image Analysis*. Springer  
619 Berlin / Heidelberg. volume 4477 of *Lecture Notes in Computer Science*,  
620 pp. 394–401.
- 621 Moeslund, T.B., Hilton, A., Krüger, V., 2006. A survey of advances in vision-  
622 based human motion capture and analysis. *Comput. Vis. Image Underst.*  
623 104, 90–126.
- 624 Oikonomopoulos, A., Patras, I., Pantic, M., 2005. Spatiotemporal salient  
625 points for visual recognition of human actions. *Systems, Man, and Cyber-*  
626 *netics, Part B: Cybernetics, IEEE Transactions on* 36, 710 –719.
- 627 Poppe, R., 2010. A survey on vision-based human action recognition. *Image*  
628 *and Vision Computing* 28, 976 – 990.
- 629 Saghafi, B., Rajan, D., 2012. Human action recognition using pose-based  
630 discriminant embedding. *Signal Processing: Image Communication* 27, 96  
631 – 111.

- 632 Schuldt, C., Laptev, I., Caputo, B., 2004. Recognizing human actions: a local  
633 svm approach, in: Pattern Recognition, 2004. ICPR 2004. Proceedings of  
634 the 17th International Conference on, pp. 32 – 36 Vol.3.
- 635 Scovanner, P., Ali, S., Shah, M., 2007. A 3-dimensional sift descriptor and its  
636 application to action recognition, in: Proceedings of the 15th international  
637 conference on Multimedia, ACM, New York, NY, USA. pp. 357–360.
- 638 Shao, L., Chen, X., 2010. Histogram of body poses and spectral regression  
639 discriminant analysis for human action categorization, in: British Machine  
640 Vision Conference (BMVC), Aberystwyth, UK.
- 641 Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore, R.,  
642 Kipman, A., Blake, A., 2011. Real-time human pose recognition in parts  
643 from single depth images, in: Computer Vision and Pattern Recognition  
644 (CVPR), 2011 IEEE Conference on, pp. 1297 –1304.
- 645 Singh, S., Velastin, S., Ragheb, H., 2010. Muhavi: A multicamera human ac-  
646 tion video dataset for the evaluation of action recognition methods, in: Ad-  
647 vanced Video and Signal Based Surveillance (AVSS), 2010 Seventh IEEE  
648 International Conference on, pp. 48 –55.
- 649 Suzuki, S., be, K., 1985. Topological structural analysis of digitized binary  
650 images by border following. Computer Vision, Graphics, and Image Pro-  
651 cessing 30, 32 – 46.
- 652 Thureau, C., Hlaváč, V., 2007.  $n$ -grams of action primitives for recognizing  
653 human behavior, in: Kropatsch, W., Kampel, M., Hanbury, A. (Eds.),

- 654 Computer Analysis of Images and Patterns. Springer Berlin / Heidelberg.  
655 volume 4673 of *Lecture Notes in Computer Science*, pp. 93–100.
- 656 Tran, D., Sorokin, A., 2008. Human activity recognition with metric learn-  
657 ing, in: Forsyth, D., Torr, P., Zisserman, A. (Eds.), *Computer Vision*  
658 *ECCV 2008*. Springer Berlin / Heidelberg. volume 5302 of *Lecture Notes*  
659 *in Computer Science*, pp. 548–561.
- 660 Turaga, P., Chellappa, R., Subrahmanian, V., Udrea, O., 2008. Machine  
661 recognition of human activities: A survey. *Circuits and Systems for Video*  
662 *Technology*, *IEEE Transactions on* 18, 1473 –1488.
- 663 Wang, J., She, M., Nahavandi, S., Kouzani, A., 2010. A review of vision-  
664 based gait recognition methods for human identification, in: *Digital Image*  
665 *Computing: Techniques and Applications (DICTA)*, 2010 International  
666 Conference on, pp. 320 –327.
- 667 Weinland, D., Özuysal, M., Fua, P., 2010. Making action recognition ro-  
668 bust to occlusions and viewpoint changes, in: Daniilidis, K., Maragos, P.,  
669 Paragios, N. (Eds.), *Computer Vision ECCV 2010*. Springer Berlin / Hei-  
670 delberg. volume 6313 of *Lecture Notes in Computer Science*, pp. 635–648.
- 671 Weinland, D., Ronfard, R., Boyer, E., 2006. Free viewpoint action recognition  
672 using motion history volumes. *Comput. Vis. Image Underst.* 104, 249–257.
- 673 Weinland, D., Ronfard, R., Boyer, E., 2011. A survey of vision-based methods  
674 for action representation, segmentation and recognition. *Comput. Vis.*  
675 *Image Underst.* 115, 224–241.

- 676 Wu, C., Khalili, A.H., Aghajan, H., 2010a. Multiview activity recogni-  
677 tion in smart homes with spatio-temporal features, in: Proceedings of  
678 the Fourth ACM/IEEE International Conference on Distributed Smart  
679 Cameras, ACM, New York, NY, USA. pp. 142–149.
- 680 Wu, X., Shi, Z., Zhong, Y., 2010b. Detailed analysis and evaluation of key-  
681 point extraction methods, in: Computer Application and System Modeling  
682 (ICCASM), 2010 International Conference on, pp. V2-562 –V2-566.
- 683 Wu, X., Xu, D., Duan, L., Luo, J., 2011. Action recognition using context  
684 and appearance distribution features, in: Computer Vision and Pattern  
685 Recognition (CVPR), 2011 IEEE Conference on, pp. 489 –496.