

1 **Common-Sense Reasoning for Human Action Recognition**

2 Jesús Martínez del Rincón[#], Maria J. Santofimia*, Jean-Christophe Nebel[‡]

3 [#]The institute of Electronics, Communications and Information Technology (ECIT),
4 Queen's University of Belfast, BT3 9DT, UK

5
6 * Department of Technology and Information Systems, Computer Engineering
7 School, University of Castilla-La Mancha, Ciudad Real, Spain

8 [‡]Digital Imaging Research Centre, Kingston University, London, KT1 2EE, UK

9 10 **Abstract**

11 This paper presents a novel method that leverages reasoning capabilities in a
12 computer vision system dedicated to human action recognition. The proposed
13 methodology is decomposed into two stages. First, a machine learning based
14 algorithm – known as bag of words- gives a first estimate of action classification from
15 video sequences, by performing an image feature analysis. Those results are
16 afterward passed to a common-sense reasoning system, which analyses, selects
17 and corrects the initial estimation yielded by the machine learning algorithm. This
18 second stage resorts to the knowledge implicit in the rationality that motivates human
19 behaviour. Experiments are performed in realistic conditions, where poor recognition
20 rates by the machine learning techniques are significantly improved by the second
21 stage in which common-sense knowledge and reasoning capabilities have been
22 leveraged. This demonstrates the value of integrating common-sense capabilities
23 into a computer vision pipeline.

24 *Keywords:* Common sense, artificial intelligence, action recognition, bag of words,
25 computer vision

27 **1. Introduction**

28 In the last decade, the automated recognition of human actions from video
29 sequences has become an essential field of research in computer vision. Not only
30 does it have applications in video surveillance, but also in indexing of film archives,
31 sports video analysis and human-computer interactions. However, the task of action
32 recognition from a single video remains extremely challenging due to the huge
33 variability in human shape, appearance, posture, the individual style in performing
34 some actions, and external contextual factors, such as camera view, perspective and
35 scene environment.

36 During the last few years, thanks to the availability of many datasets suitable for
37 training action recognition algorithms, the field has made enormous progress to the
38 point that the automatic annotation of the KTH (Schuldt et al., 2004) and Weizzman
39 (Blank et al., 2005) databases is now considered solved. For more complex data, i.e.
40 IXMAS (Weinland et al., 2006) and UT-Interaction (Ryoo and Aggarwal, 2009),
41 accuracy rates around 80% are now claimed by state-of-the-art approaches
42 (Waltisberg et al., 2010; Weinland et al., 2010; Nebel et al., 2011). Unfortunately, all
43 those action recognition experiments are conducted with videos that are not
44 representative of real life data, which led a recent review to conclude that none of
45 existing techniques would be currently suitable for real visual surveillance
46 applications (Nebel et al, 2011). This is further confirmed by the poor performance,
47 obtained on videos captured in uncontrolled environments, such as Hollywood 1 and
48 2 datasets (Laptev et al. 2008) and Human Motion DataBase (HMDB51) (Kuehne et
49 al., 2011), where accuracies are 32%, 51% and 20% respectively (Kuehne et al.,

50 2011). In addition, these challenging datasets only display a fraction of the
51 complexity exhibited by the real world, e.g. at most 51 different actions are
52 considered. Consequently, usage of video-based action recognition remains a very
53 distant aspiration for most actual applications.

54 On the other hand, the human brain seems to have perfected the ability to recognise
55 human actions despite their high variability. This capability relies not only on
56 acquired knowledge, but also on the aptitude of extracting information relevant to a
57 given context and logical reasoning. In contrast, machine learning based action
58 recognition methodologies tend to learn isolated actions from a set of examples.
59 Although only a few and limited attempts to introduce contextual information have
60 been made (Waltisberg et al., 2010; Chen and Nugent, 2009; Akdemir et al. 2008;
61 Vu et al. 2002; Ivano and Bobick, 2000), their performance supports the idea that
62 action recognition can benefit greatly from combining traditional computer vision
63 based algorithms with knowledge based approaches.

64 In this paper, we propose a novel method relying on common-sense reasoning and
65 contextual and common-sense knowledge which allows analysing, selecting and
66 correcting annotation predictions made by a video-based action recognition
67 framework. The presented approach is decomposed into two stages. First, a classic
68 action recognition algorithm classifies actions independently according to similarity to
69 the training set. Secondly, results are refined using common-sense knowledge and
70 reasoning. More specifically, contextual information is exploited using common
71 sense reasoning.

72 **2. Relevant work**

73

74

a. Video-based Human Action Recognition

75 Video-based activity recognition algorithms can be classified into two different
76 classes: those that train from examples and those that provide descriptions of
77 general types . The first and main category includes action descriptors based on
78 Hidden Markov Models (Vezzani et al., 2010; Kellokumpu et al, 2008; Martinez et al.
79 2009; Ahmad and Lee, 2008; Weinland et al., 2007), Conditional Random Field
80 (Zhang and Gong, 2010; Natarajan and Nevatia, 2008; Wang and Suter, 2007), Bag
81 of Words (Laptev et al., 2008; Liu and Shah, 2008; Matikainen et al., 2010; Ta et al.,
82 2010; Liu et al., 2008; Kovashka and Grauman, 2010) and low dimension manifolds
83 (Wang and Suter, 2007b, 2008; Fang et al. 2009; Jia and Yeung, 2008; Blackburn
84 and Ribeiro, 2007; Richard and Kyle, 2009; Turaga et al. 2008; Lewandowski et al.
85 2010, 2011). Since those approaches do not include any reasoning capability, their
86 efficiency relies on a training set which is supposed to cover the variability of all
87 actions present in the target videos. Given that this condition can only be valid in the
88 most controlled scenarios, it has been proposed to extend these techniques by
89 adding some form of reasoning based on either rules or logic.

90 The inclusion of reasoning has been sparsely used and mostly for specific
91 applications. It should be noted it is particularly popular in intelligent surveillance for
92 the detection of unusual events (Makris et al. 2008). Since training data do not exist
93 to define those events, rules and reasoning are the only available tools. Usually,
94 activities which do not match those present in the training set are classified as
95 unusual. In the most specific field of action recognition, reasoning rules have proved
96 particularly successful when dealing with interactions between subjects (Waltisberg
97 et al. 2010). Indeed, following initial action recognition on each character individually
98 using a Random Forest framework, analysis of those actions allows inferring the

99 nature of their interaction. As reported by Waltisberg et al. (2010), this scheme
100 outperforms the standard approach which deals with all characters at once and is the
101 current state of the art on the UT-Interaction dataset (Ryoo and Aggarwal, 2009).
102 These results support our hypothesis that additional knowledge and reasoning will
103 lead to better performance.

104 The second class of video-based activity recognition algorithms exploits a common
105 knowledge-base or ontology of human activities to perform logical reasoning. Since
106 ontology design is empirical in nature and labour intensive - symbolic action
107 definitions are based on manual specification of a set of rules -, current ontologies
108 are only suitable for very specific scenarios. In the field of video surveillance,
109 ontologies have been proposed for analysis of social interaction in nursing homes
110 (Chen et al., 2004), classification of meeting videos (Hakeem and Shah, 2004) and
111 recognition of activities occurring in a bank (Georis et al., 2004). However, there is a
112 need for an explicit commonly agreed representation of activity definitions
113 independently of domain and/or algorithmic choice. Such common knowledge base
114 and its exploitation through rules would facilitate portability, interoperability and
115 sharing of reasoning methodologies applied to activity recognition. Several attempts
116 have been made to design ontologies for visual activity recognition in a more
117 systematic manner (Akdemir et al., 2008, Hobbs et al., 2004, Francois et al, 2005) so
118 that they can cover different scenarios, e.g. both bank and car park monitoring
119 (Akdemir et al., 2008). However, they remain limited to a few domains - up to 6
120 (Hobbs et al., 2004).

121

122 **b. Common Sense Reasoning**

123 Within the artificial intelligence (AI) community, the usage of video as information
124 source for reasoning has not been extensively applied (Moore et al., 1999; Duong et
125 al., 2005). This is due to the lack of robustness and consistency of video features in
126 real world scenarios, where the huge variability of the conditions impact considerably
127 on activity recognition. As a consequence, AI researchers have focused on using
128 sensors which are more reliable and consistent, but more intrusive, sensors to
129 gather an actor's behavioural information (Wang et al. 2007c). They include
130 wearable sensors based on inertial measurement units (e.g. accelerometers,
131 gyroscopes, magnetometers) and RFID tags attached to the actors and/or to objects.
132 In such set-up, complex reasoning is possible and successful artificial intelligence
133 approaches have flourished (Wang et al., 2007c; Philipose et al., 2004; Tapia et al.,
134 2004). However, most of these sensors are not suitable in most real life applications
135 due to either their intrusive nature, e.g. subjects may refuse to wear them, or
136 technical factors, such as size, ease of use and battery life.

137 Among the AI approaches which could be considered for video based human action
138 recognition, common-sense, probabilistic and ontological reasoning, as described in
139 the previous subsection, are of particular interest. Ontological languages such as
140 OWL (Dean et al., 2011a) and RDF (Dean et al., 2011b) use a syntax that imposes
141 severe restrictions in the type of information that can be represented. First,
142 relationships involving more than two entities cannot be considered since they may
143 lead to hold a-priori inconsistent information, which is not allowed in this
144 methodology. Secondly, since reasoning is limited to checking the consistency of the
145 knowledge base, new information cannot be inferred. Both common-sense and
146 probabilistic reasoning are able to address those limitations. However, their nature is
147 very different since they can be classified as techniques based on either qualitative

148 or quantitative reasoning. A weakness of quantitative reasoning comes from the
149 complexity of estimating accurate probabilities for activities of interest: in practice it is
150 unfeasible when dealing with unconstrained and realistic scenarios (Kuipers, 1994).
151 On the other hand, qualitative reasoning has the ability of considering causality and
152 expected behaviour based on logics, i.e. reasoning can provide explanations
153 rationalising or motivating a given action, whereas probabilistic reason can only
154 support decisions according to probability associated to actions.

155 As a consequence, common-sense reasoning (McCarthy, 1968, 1979; Minsky, 1986;
156 Lenat, 1989, 1990) appears particularly suited to video based human action
157 recognition. It provides the capability of understanding the context situation, given
158 the general knowledge that dictates how the world works, which allows correcting
159 mistakes made by the video analysis system. McCarthy proposes an approach to
160 build a system with the capability to solve problems in the form of an “advice taker”
161 (McCarthy, 1968). In order to do so, he reckons that such an attempt should be
162 founded in the knowledge of the logical consequences of anything that could be told,
163 as well as the knowledge that precedes it. In that work, he postulates that “a program
164 has common sense if it automatically deduces from itself a sufficiently wide class of
165 immediate consequences of anything it is told and what it already knows”. Following
166 McCarthy and Minsky’s studies (McCarthy, 1968; Minsky, 1986), it appears a way of
167 enhancing systems with the capability to understand and reason about the context is
168 by introducing commonsense knowledge similar to that humans hold.

169 In this work, we propose the integration of common-sense knowledge and reasoning
170 within a video human activity recognition framework in order to improve accuracy.
171 First, a machine learning based action recognition algorithm processes videos to
172 generate data appropriate for logical inferences. Consequently, video data become a

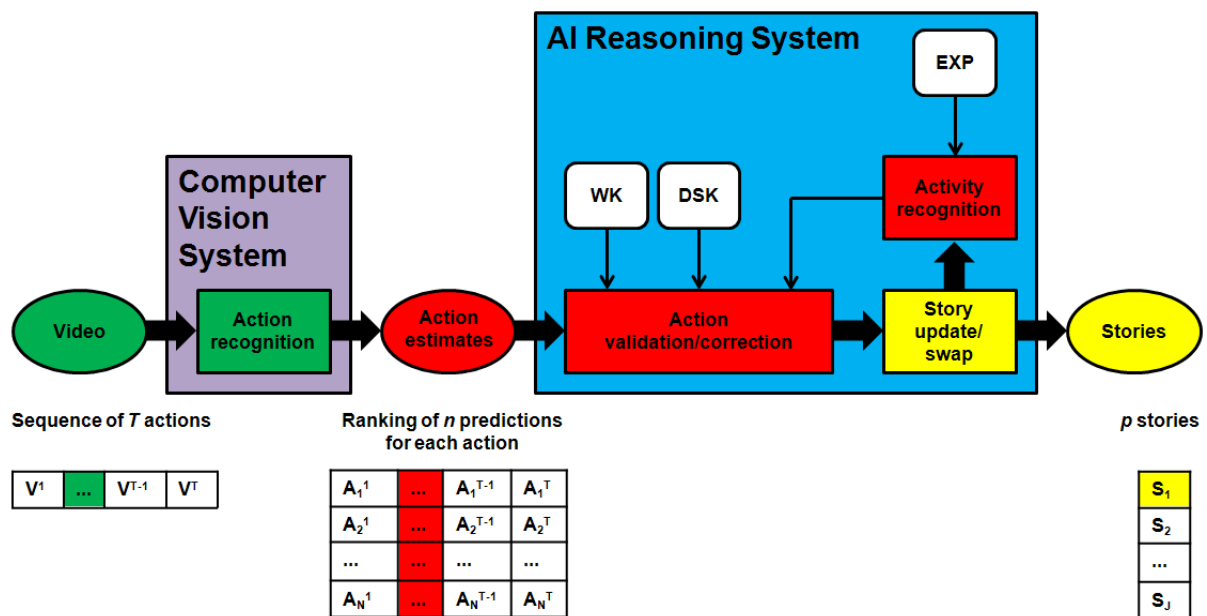
173 suitable information source for reasoning. Secondly, common-sense reasoning
 174 increases accuracy of the computer vision algorithm by introducing general, so
 175 called common-sense, and context-independent knowledge. This addition should
 176 allow usage of video based systems within real life applications.

177 3. Novel action recognition framework

178

179 a. Principles

180 We propose a novel two-stage framework where initial action predictions made by a
 181 machine learning approach are analysed, refined and, possibly, corrected by the
 182 second layer common-sense reasoning system.



183

184

Figure 1: Action recognition framework

185 Given a video, V , which can be divided into a sequence of T actions and a computer
 186 vision system (CVS) trained to recognise N types of actions, each action, V^t , is
 187 processed independently and is associated to an action estimation vector, A^t , which
 188 ranks the N types of actions according to their similarity to V^t . Eventually, the CVS

189 generates an action estimation matrix, A , of dimensions $(T \times N)$, where A_j^t represents
190 the j^{th} most likely type of the t^{th} action occurring in the video. Each action estimate
191 generated by the CVS is passed as input to the AI reasoning system (AIRS) which
192 produces, in an online manner, J stories, S_j . These stories are generated and
193 updated according to every new estimate A^t .

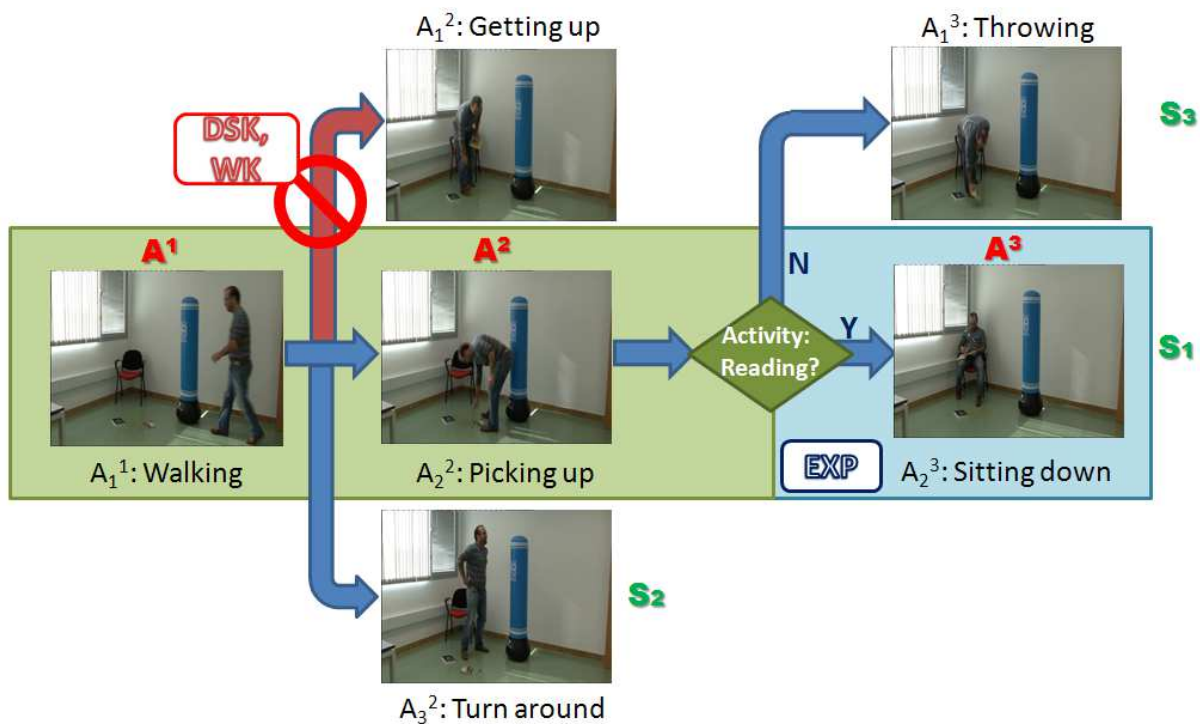
194 In this paper, we define a 'story' as a coherent list of action types describing a video
195 of interest. Coherence is defined by respect to both world and domain specific
196 knowledge, WK and DSK respectively. Selection of action types relies on common-
197 sense reasoning applied to the action estimations A , and possible recognition of
198 activities defined in the expectation knowledge, EXP. Note that a story may contain
199 'unknown action' labels when, for a given action, none of the estimations allows
200 coherent annotation. Stories are ordered by the AIRS and the most likely one is
201 always first, in the same way that actions have been ordered and prioritised by the
202 CVS.

203 The AIRS processes every action estimation vector, A^t , according to the J stories S_j
204 existing at $t-1$. First, the validity of each action estimates A_j^t is verified within the
205 context of each story S_j using knowledge contained in WK and DSK. This is done
206 inside the block Action validation/correction depicted in Figure 1. Secondly, if the
207 sequence of previous actions stored in S_j led to the recognition by EXP of an activity
208 (Figure 1, block Activity Recognition) which expected a specific action type in order
209 to be completed, and if that type is not present in A^t , a correction of A^t is performed,
210 i.e. the expected type is added to the story S_j instead of A^t . Finally, each valid action
211 of A^t updates an existing story (Figure 1, block story update/swap). If a valid action
212 cannot be allocated to a story, a new story is created. Since during the process, the
213 most likely action estimates have priority to be allocated to the first stories, S_1 is the

214 story which is the most likely to describe accurately the video of interest. However, if
 215 any other S_j shows a more likely storyline, the position of S_1 as 'main story' may be
 216 swapped with S_j (Figure 1, block story update/swap).

217 We illustrate some of the reasoning performed by AIRS with an example, see Figure
 218 2: an activity ('Getting up') incompatible with the current story (S_1) is rejected
 219 according to the world and domain specific knowledge; valid actions ('Throwing' &
 220 'Sitting down') are assigned to parallel stories (S_2 and S_3); an activity ('Reading') is
 221 recognised based on expectations, consequently the expected action ('Sitting down')
 222 is prioritised.

223



224

225 Figure 2: Example of reasoning performed by AIRS. Blue and red arrows represent,
 226 respectively, valid and invalid actions. Green box depicts the sequence of action
 227 which led to the recognition of an activity (reading) based on expectations. Blue box
 228 shows the expected action (sitting down).

229

b. Common sense reasoning algorithm

230 The AIRS assigns and evaluates correspondences between action estimations in
231 vector A^t and the stories S existing at $t-1$. The validity of each action estimate A_i^t is
232 verified sequentially within the context of the main story S_1 using knowledge
233 contained in WK and DSK. Once action allocation, if any, has been completed for the
234 main story, the same process is followed for all the other stories S_j using the
235 remaining action estimates. This double sequentiality in the assignment of actions to
236 stories deals with the fact that both stories and actions are ordered, where the first
237 actions/stories are always the most likely.

238 The n first action estimates are all considered as possible alternatives. Therefore,
239 new stories are created if they do not fit any of the existing ones. The rationale
240 behind this is that, although the first estimate provided by the CVS is not always
241 correct, the CVS is quite robust since the correct action is likely to be present among
242 the first n estimates (see 'Experimental results' section). During the allocation
243 process of a given time step, some stories may not be allocated to any action, if
244 none of the available action estimates is valid in their context according to WK and
245 DSK.

246 A second level of reasoning is introduced by exploiting the concept of activity
247 recognition. This is modelled in our system through the expectation knowledge, EXP.
248 For each story S_j , if the sequence of previous actions leads to the recognition of an
249 activity by EXP, the next action assigned to the story S_j must match the expected
250 one, eA . In case where the expected action type is not present in A^t , A^t is corrected
251 by including eA in the estimate vector so that eA can be assigned to story S_j . This
252 mechanism provides a higher level of reasoning, going further than the validation
253 mechanism provided by the DSK and WK, which allows correcting estimate failures
254 of the CVS. However, in order to avoid over-reasoning errors, corrections are

255 introduced only when, in addition to validation, a unique activity is recognised, i.e.
256 when there is no doubt regarding the type of the expected action.

257

258 Through the previously described process, the AIRS gives priority to the most likely
259 action estimates in their allocations to the first stories. As a consequence, the AIRS
260 output is an ordered set of stories, where S_1 is the story which is the most likely to
261 describe accurately the video of interest.

262 However, the accuracy of the CVS may depend of the nature of the action and vary
263 over time during video processing, which may lead to the correct estimates to be
264 lower in the action estimation vectors. Consequently, after a while S_1 may not
265 contain the most likely story. The AIRS addresses this issue using a story swapping
266 mechanism. When the AIRS is able to allocate systematically actions to a given story
267 S_j and activities kept being recognised according to the expectations, this story is
268 accepted as the main story and swapped with S_1 . Empirical experimentations have
269 shown that the story swapping mechanism should be triggered when a story displays
270 two consecutive activity recognitions, $TH=2$.

271

272 This reasoning algorithm is presented through the following pseudo code. First, the
273 main variables are defined. Then, the core of the algorithm is detailed. Finally, the
274 main functions are described. Note that functions are colour-coded to allow better
275 readability of the algorithm.

```
276  
277 //////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////  
278 // INPUT  
279 //////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////  
280 // Expert systems  
281 Expert DSK,WK,Exp;  
282 //An action is a primitive  
283 Action eA; // expected action  
284 Action At[N]; // alternative actions predicted for time t,  
285 // At are ranked according to CVS's prediction confidence
```

```

286 Int N;           // number of alternative actions at time t
287 //A story is a list of actions
288 Story S[J];     // existing stories
289 Int J=1;       // number of existing stories, one starts with 1 story
290 S[1]=null;     // the initial story is empty
291
292 //Each story is associated to a list of possible activities containing
293 future actions for the next time t
294 Typedef Action[] Activity;
295 Activity PossibleActiv[][J]=[ ALL ][J]; // set of activities, initially all
296                                     // activities are possible
297 Int expect_fulfill[J]=zeros(1,J); // story counter for swapping mechanism
298 ////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
299 // MAIN
300 ////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
301 for t=1:Inf           // for each time step
302     N=length(At);           // number of alternative actions
303     Bool assigned_action[N]=zeros(1,N); // no action is assigned
304     J=length(S);     // number of existing stories
305     Bool updated_story[J]=zeros(1,J); // no story has been updated
306     for i=1:N        // for each alternative action
307         // integration of action i into an existing story
308         for j=1:J    // for each existing story
309             if (updated_story(j)==0) // if story j is available
310                 // activity recognition process
311                 eA=f_activity_recognition(PossibleActiv(j)); //expected activity
312                 if (eA!=null) // if activity recognised //
313                     story updating process
314                     [PossibleActiv(j),S(j)]= f_story_update
315                                     (eA,PossibleActiv(j),S(j),Exp);
316                     updated_story(j)=1; // story j is updated
317                     // action allocation process
318                     assigned_action=f_action_allocation(assigned_action,eA,At);
319                     // story swapping process
320                     [S,expect_fulfill]=f_storySwapping(S,expect_fulfill,j);
321                 else // no activity is recognised
322                     if (assign_action(i)==0) // if action i is available
323                         // action validation process
324                         if f_action_validation(At(i),DSK,WK,S(j)) //if At(i) valid
325                             // story updating process
326                             [PossibleActiv(j),S(j)]=f_story_update
327                                     (At(i),PossibleActiv(j),S(j),Exp);
328                             updated_story(j)=1; // story j is updated
329                             // action allocation process
330                             assign_action(i)=1; // action i is allocated
331                         end
332                     end
333                 end
334             end
335         end
336         // integration of non-assigned action i into a new story
337         if (assign_action(i)==0) // if action i is available
338             // action validation process
339             if f_action_validation(At(i),DSK,WK,S(j)) // if action i is valid
340                 // story creation process
341                 [PossibleActiv,S,expect_fulfill]=f_story_creation
342                                     (S,At(i),Exp,expect_fulfill);
343                 J=length(S); // update number of stories
344                 updated_story(J)=1; // story J is updated
345                 // action allocation process
346                 assign_action(i)=1; // action i is allocated

```

```

347         end
348     end
349 end
350 end

```

351 Expectations are checked at each given time t , for each current story (function
352 `f_activity_recognition`). If the number of current expected activities is only one,
353 the nature of the ongoing activity is known. Therefore, the function is able to return
354 the expected type of the next action, eA .

```

355 function [Action a]=f_activity_recognition(Activity pred)
356     if (size(pred)==1)
357         return pred(1);
358     else
359         return null;
360     end

```

361 If any of the n observed actions of A^t matches eA , this action is set as allocated to
362 avoid inclusion in any other story (function `f_action_allocation`).

```

363 function [bool b]=f_action_allocation(bool b, Action a, Action[] v)
364     for i=1:size(v)
365         if(v(i)==a)
366             b=1;
367         end
368     end
369     return b;

```

370 When an action has been judged suitable to be added to a story, the current story is
371 updated (function `f_story_update`). This also involves updating the list of possible
372 ongoing activities, i.e. knowledge about possible actions for time $t+1$:
373 `PossibleActiv(j)`. This is achieved by, first, retrieving all expected activities in the
374 knowledge of action a at time t , $p2$, (function `retrieve_expected_activities`)
375 and, then, by finding the intersection between this list and the one predicted for time
376 t , p , (function `intersection`). If no intersection exists, i.e. either CVS has failed or
377 reasoning has been erroneous, since it is not possible to distinguish the source of
378 the failure, expected activities are reset to $p2$ to avoid propagating errors.

```

379 function [Activity p,Story s]=f_story_update
380     (Action a, Activity p, Story s, Exp e)

```

```

381     Activity p2=null;
382     s=[s a];           // add action a to current story s
383     p2=retrieve_expected_activities(e,a);
384     p=intersection(p,p2); // new list of expected activities
385     if (size(p)==0)
386         p=p2;
387     end;
388     return [p,s];

```

389 If the activity recognition algorithm was able to detect unequivocally the nature of an
390 ongoing activity within a story, S_j , confidence in that story is increased. This is stored
391 in the variable `expect_fulfill`. The value of that variable is evaluated during the
392 story swapping mechanism (function `f_storySwapping`). If it shows that the story S_j
393 has consecutively recognised activities (in our case twice $TH=2$), the story S_j is
394 swapped with S_1 and becomes the main story, i.e. the most likely one.

```

395 function [Story s[], int[] f]=f_storySwapping(Story s[], int[] f, int indx)
396     Story s_tmp;
397     f(indx)++;
398     if f(indx)>=TH
399         // s(indx) is moved as top story and all the others are shifted down
400         s = [s(indx) s(1: indx-1) s(indx-1:end)];
401         f = zeros(1,J);
402     end
403     return [s,f];

```

404 If the activity recognition mechanism does not detect any ongoing activity or several
405 activities are possible, action allocation only relies on action validity. This is
406 evaluated according to the action global coherence with the world WK and the
407 domain specific knowledge DSK within the context of a story (function
408 `f_action_validation`).

```

409 function bool=f_action_validation(Action a,DSK d,WK w,Story s)
410     return validate(a,d,s,w);

```

411 If an action is judged as valid, the action is assigned to the story and expected
412 activities are updated (function `f_story_update`). After the assignment, boolean
413 vectors, `assigned_action` and `updated_story`, are updated to make sure that each
414 action is assigned at most to one story and that each story is not updated more than
415 once for a given time t .

416 Finally, if an action is valid but has not been assigned to any current story, a new
417 story is created (function `f_story_creation`).

```
418 function [Activity p, Stories s, int[] f]=f_story_creation(Stories s,  
419 Action a, EXP e, Activity p, int[] f)  
420     Activity Activnew=[All];  
421     Story Snew=[];  
422     [Activnew, Snew]=f_story_update(a,Activnew,Snew,e);  
423     J=J+1;  
424     s(J)=Snew;  
425     p(J)= Activnew;  
426     expect_fulfill(J)= 0;  
427     return [p,s];
```

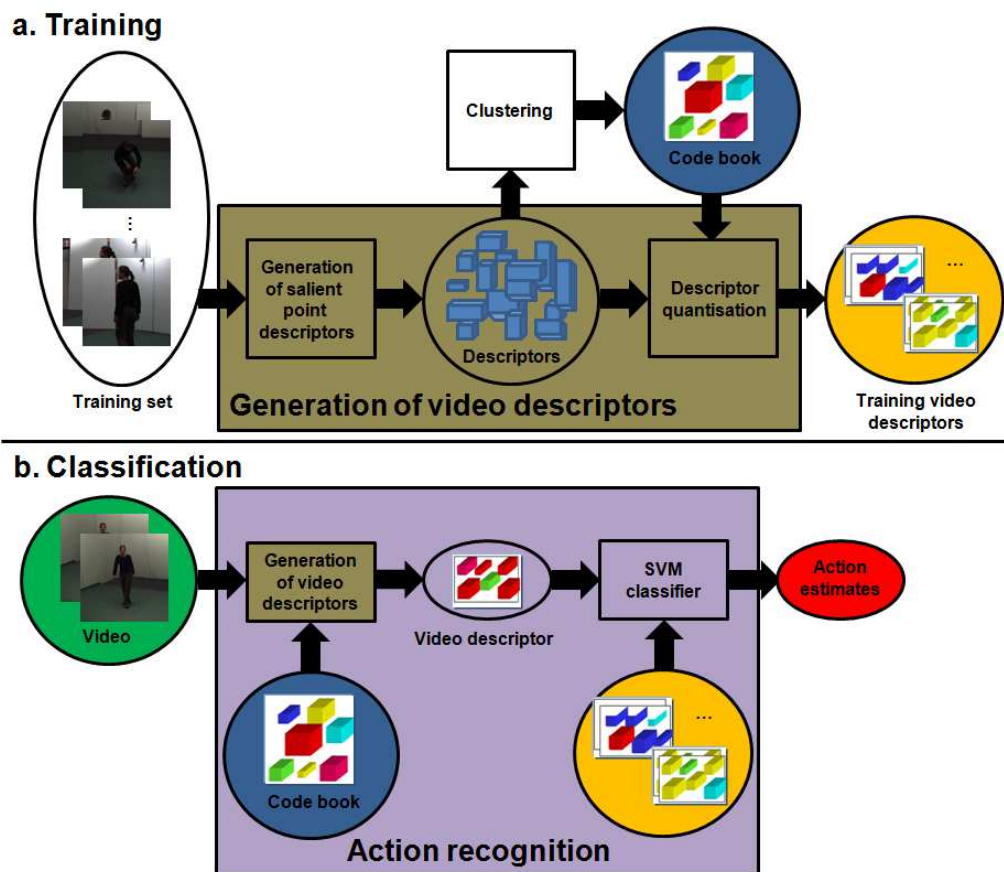
428 **4. Implementation**

429

430 **a. Computer vision based action recognition**

431 Although computer vision based action recognition has been a very active field of
432 research, only a few approaches have been evaluated on view independent
433 scenarios. Accurate recognition has been achieved using multi-view data with either
434 3D exemplar-based HMMs (Weinland et al., 2007) or 4D action feature models (Yan
435 et al. 2008). But, in both cases performance dropped significantly in a monocular
436 setup. This was addressed successfully by representing videos using self-similarity
437 based descriptors (Junejo et al., 2008). However, this technique assumes a rough
438 localisation of the individual of interest which is unrealistic in many applications.
439 Similarly, the good performance of a SOM based approach using motion history
440 images is tempered by the requirement of segmenting characters individually (Orrite
441 et al. 2008). More recently a few approaches have produced accurate action
442 recognition from simple extracted features: two of them rely on a classifier trained on
443 bags of words (Kaaniche and Bremond, 2010; Liu et al. 2008) whereas the other one
444 is based on a nonlinear dimensionality reduction method designed for time series
445 (Lewandoski et al. 2010).

446 Among those approaches, the Bag of Words (BoW) framework is particularly
 447 attractive since, not only it is one of the most accurate methods for action
 448 recognition, but its computational cost is low. Moreover, BoW can be applied directly
 449 on video data without the need of any type of segmentation. The versatility of that
 450 framework has been demonstrated on a large variety of datasets including film-
 451 based ones (Laptev and Perez, 2007). Consequently, in this study, we decided to
 452 base the computer vision system of our action recognition framework on a BW
 453 methodology.



454
 455 Figure 3: BoW framework: a) Training and b) classification pipelines

456 BoW is a learning method which was used initially for text classification (Joachims,
 457 1998). It relies on, first, extracting salient features from a training dataset of labelled
 458 data. Then, these features are quantised to generate a code book which provides

459 the vocabulary in which data can be described and analysed. Here, we based our
460 implementation on that proposed by (Csurka et al., 2004).

461 The BoW training stage aims at, first, producing a codebook of feature descriptors
462 and, secondly, generating a descriptor for each action video available in the training
463 set, see Figure 3 a). The training pipeline starts by detecting salient feature points in
464 each video using a spatio-temporal detector (Harris 3D) and describing each
465 individual point by a histogram of optical flow (STIP) (Laptev, 2005). Once feature
466 points are extracted from all training videos, the k-means algorithm is employed to
467 cluster the salient point descriptors into k groups, where their centres are chosen as
468 group representatives. These points define the codebook which is then used to
469 describe each video of the training set. Finally, those video descriptors are used to
470 train SVM classifiers – one per action of interest - with a linear kernel.

471 In order to recognise the action performed in a video, Figure 3 b), salient feature
472 points are first detected. Then, their descriptors are quantified using the codebook in
473 order to generate a video descriptor. Finally, the video descriptor is fed into each
474 SVM classifier, which allows quantifying the fit between the video and each trained
475 action type. Therefore, an action estimation vector A can be generated where action
476 types are ranked according to their fit.

477 **b. Knowledge-Base System for Common Sense Reasoning**

478 Automating common-sense reasoning requires an expressive-enough language, a
479 knowledge base and a set of mechanisms capable of processing this knowledge to
480 check consistency and infer new information. A few knowledge-based approaches
481 offer such features, i.e. Scone (Chen and Fahlman, 2008; Fahlman, 2006), Cyc
482 (Lenat et al. 1989, 1990), WordNet (Fellbaum, 1998) or ConceptNet (Eagle et al.,

483 2003). Among them, the open-source Scone project is of particular interest since,
484 instead of placing its focus on collecting common-sense knowledge, it provides
485 efficient and advanced means for accomplishing search and inference operations.

486 The main difference between this and other approaches lies in the way in which
487 search and inference are implemented. Scone adopts a marker-passing algorithm
488 (Fahlman, 2006), which is not a general theorem-prover, but is much faster and
489 supports most of the search and inference operations required in common-sense
490 reasoning: inheritance of properties, roles, and relations in a multiple-inheritance
491 type hierarchy; default reasoning with exceptions; detecting type violations; search
492 based on set intersection; and maintaining multiple, immediately overlapping world-
493 views in the same knowledge base. In addition, Scone provides a multiple-context
494 mechanism which emulates humans' ability to store and retrieve pieces of
495 knowledge, along with matching and adjusting existing knowledge to similar
496 situations.

497 In our framework, the algorithm described in section 3b was implemented using
498 Scone in order to encode formal definitions and their applications for WK, DSK and
499 EXP. It is important to note that, although we took advantage of the proposed multi-
500 context mechanism (Chen and Fahlman, 2008), we exploited it for a usage it was not
501 originally intended for, extending its application for a wider purpose. In particular, we
502 propose the usage of multi-context for the management of alternative stories
503 describing coherent explanations of the video of interest.

504 The three sources of knowledge exploited in our implementation, i.e. WK, DSK and
505 EXP, are described below:

506 1. World knowledge, WK, comprises all relevant common-sense knowledge that
507 describes “how the world works”. This information is independent of the
508 application domain, in the sense that it only considers general knowledge
509 rather than specific or expert knowledge about a specific field. As an example,
510 we provide below the description of the implications of performing the action
511 of ‘scratching the head’.

```
512 (new-event-type {scratch} '({event}))  
513 :roles  
514 ((:type {scratcher} {animated thing})  
515 (:type {scratched thing} {thing})))  
516 (new-event-type {scratch head}  
517 '({scratch} {action}))  
518 :roles  
519 ((:rename {scratched thing} {scratched head})  
520 (:rename {scratcher} {scratcher hand}))  
521 :throughout  
522 ((new-is-a {scratcher hand} {hand}))  
523 :before  
524 ((new-statement {scratcher hand} {approaches} {scratched head})  
525 (new-not-statement {scratcher hand} {is in direct contact to}  
526 {scratched head}))  
527 :after  
528 ((new-statement {scratcher hand} {is in direct contact to}  
529 {scratched head})))
```

530 2. Domain specific knowledge, DSK, describes a given application domain in
531 terms of the entities that are relevant for that specific context, as well as, the
532 relationships established among those. The description of an element
533 “punching ball” as part of the layout of a specific room is an example of
534 domain specific information.

```
535 (new-type {bouncing element} {thing})  
536 (new-type {punching ball} {thing})  
537 (new-is-a {punching ball} {bouncing element})  
538 (new-indv-role {punching ball location} {punching ball} {location})  
539 (new-statement {punching ball} {is in} {test room})  
540 (new-statement {punching ball} {rests upon} {test room floor})  
541
```

542 3. Expectations, EXP, consist in sequences of actions that are expected to
543 happen one after the other. It encapsulates logical concepts such as causality,
544 motivation and rationality, which are expected in human action recognition.

545 For example, in a waiting room context, if a person picks up a magazine, that
546 person is expected to sit down and read the magazine. Expectations are part
547 of the domain specific knowledge since described behavioural patterns are
548 context specific.

```
549 (new-indv {picking up a book for reading it} {expectations})  
550 (the-x-of-y-is-z {has expectation} {picking up a book for reading it} {walk  
551 towards})  
552 (the-x-of-y-is-z {has expectation} {picking up a book for reading it} {pick  
553 up})  
554 (the-x-of-y-is-z {has expectation} {picking up a book for reading it} {turn  
555 around})  
556 (the-x-of-y-is-z {has expectation} {picking up a book for reading it} {sit  
557 down})  
558 (the-x-of-y-is-z {has expectation} {picking up a book for reading it} {get  
559 up})  
560
```

561 **5. Experimental results**

562

563 **i. Dataset and Experimental Setup**

564 In order to perform action recognition experiments which are relevant to real life
565 applications, videos under study should display realistic scenarios. In addition, a
566 suitable training set must be available, i.e. it must be able to cover a variety of
567 camera views so that recognition is view-independent and the set should include a
568 sufficiently large amount of instances of the actions of interest. These instances must
569 be not only annotated but perfectly segmented and organised to simplify the training.

570 The only suitable training sets which fulfil these requirements are IXMAS (Weinland
571 et al., 2006) and Hollywood (Laptev et al. 2008), as stated in the introduction.
572 Whereas the Hollywood dataset is oriented towards event detection which includes
573 significant actions but largely independent from each other (drive car, eat, kiss,
574 run...), IXMAS is focused on standard indoor actions which allows providing quite an
575 exhaustive description of possible actions in this limited scenario. Therefore, IXMAS

576 actions may be combined to describe simple activities, i.e. sit down-get up, pick up-
577 throw, punch-kick and walk-turn around, and eventually provide complete
578 representations of sets of actions performed by individual, i.e. recognition of whole
579 stories.

580 Thus, for training, the publicly available multi-view IXMAS dataset is chosen
581 (Weinland et al., 2006). It is comprised of 13 actions, performed by 12 different
582 actors. Each activity instance was recorded simultaneously by 5 different cameras.

583 Since no suitable standard videos are available in order to describe the complexity of
584 a real life application with a significant number of complex activities, we create a new
585 dataset, called the Waiting Room dataset “WaRo11” (Santofimia et al., 2012), that
586 we make available to the scientific community. In addition, using very different
587 datasets for training and testing allows us to show the generality of our framework,
588 its capabilities for real-world applications and its performance under a challenging
589 situation.

590 Since the “WaRo11” dataset has been designed for being representative of the
591 variability existing in a real life scenario, but also for integrating most of the actions
592 trained for the CVS, a specific setup was configured to simulate a waiting room. In
593 this setup, actions happen without giving any instructions to the subjects. They are
594 performed as natural part of their behaviour and motivation as human beings. This is
595 facilitated thanks to the presence of several elements interrelated to each other,
596 which may introduce causality and sequentiality as it is found in a real situation. For
597 instance, the presence of a book and a chair could motivate a subject to first pick up
598 the book and then sit down to carry out the action reading. Alternatively, a subject
599 may pick up the book, realises its topic of no interest and decides to throw it away.

600 This waiting room setup was implemented in a single room and filmed by a single
 601 fixed camera. A book was positioned on the floor, a chair was left in a corner and a
 602 punching ball was placed in another corner. Eleven sequences were recorded with
 603 eleven different actors of both genders comprising a wide range of ages (19-57) and
 604 morphological differences. No instruction was given to the actors further than “go to
 605 the room and wait for 5 minutes and feel free to enjoy the facilities while you wait”.
 606 The resulting variability in the actions performed is depicted in Table 1.

Sequence	Age	Sex	Number of actions
Actor 1	34	M	31
Actor 2	33	M	25
Actor 3	35	M	10
Actor 4	57	F	12
Actor 5	19	M	9
Actor 6	19	M	18
Actor 7	20	F	15
Actor 8	19	M	9
Actor 9	22	F	5
Actor 10	19	M	12
Actor 11	20	F	9
Total			155

Actions	Instances
check watch	4
cross arms	0
scratch head	2
sit down	13
get up	12
turn around	18
walk	53
wave hand	9
punch	26
kick	10
point	3
pick up	13
throw	0

607 Table 1: a) Number of actions performed by each actor. b) Number of instances of
 608 the trained actions found in the WaRo11 dataset.

609 Each of the recorded sequence was manually groundtruthed: first, the video of
 610 interest was segmented into a set of independent actions, then each action was
 611 labelled. Note that the segmentation of a video into independent actions is outside
 612 the scope of this study. Therefore, when testing our algorithms, we processed
 613 manually segmented actions. Readers interested in automatic action segmentation
 614 should refer to (Rui and Anandan, 2002; Black et al., 1997; Ali and Aggarwal, 2001;
 615 Shimosaka, 2007; Shi, 2011).

616 ii. Results

617 a) Performance of the computer vision system

619 First the CVS was applied to IXMAS sequences using the leave-one-out strategy
 620 followed by (Weinland et al., 2007; Yan et al., 2008; Junejo et al., 2008; Richard and
 621 Kyle, 2009). In each run, we select one actor for testing and all remaining subjects
 622 for training. Secondly, using the whole of the IXMAS dataset for training, the CVS

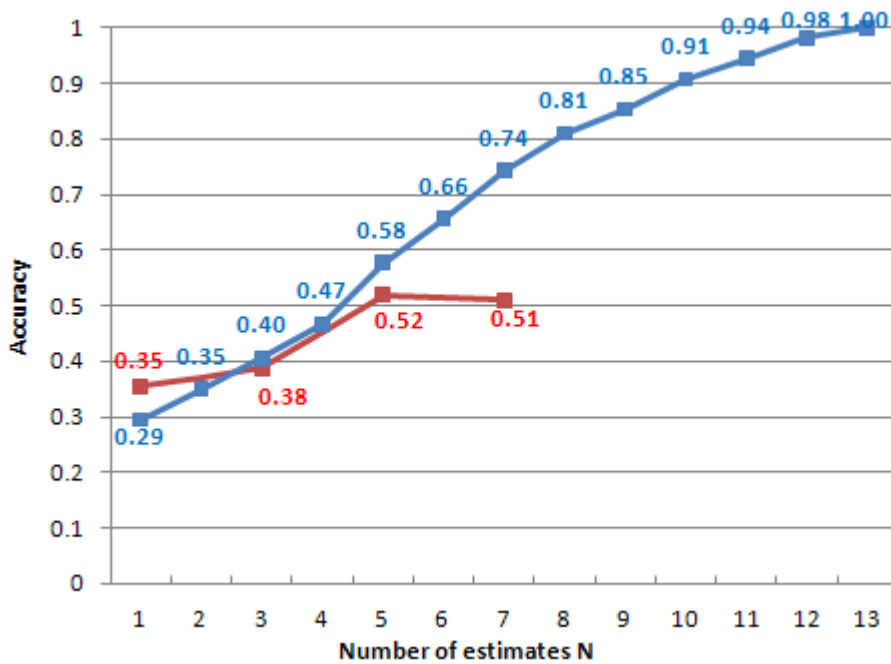
623 was applied to WaRo11. Accuracy performances for both experiments are provided
624 in Table 2.

625 Table 2. Average recognition rate for all the actions on the datasets obtained by the
626 computer vision system based on BoW

	IXMAS	WaRo11
CVS: BoW	63.9%	29.4%

627

628 The BoW based technique displays results comparable to those of the state of the
629 art on the IXMAS dataset (Nebel et al. 2011). However, when applied to a more
630 realistic environment, performances decrease considerably. This shows the
631 limitations of the CVS methodology under real circumstances, when the testing
632 conditions differs significantly from the training. On the other hand, when
633 performance is analysed in terms of average cumulative recognition curve (ACR) -
634 Figure 4, blue -, i.e. percentage that an action is accurately recognised within a set of
635 estimates,- one can see that considering the first few ranks may improve significantly
636 accuracy. For example, accuracy would jump from 29 to 66% if the best solution
637 could be detected within the 6 first estimates. This confirms that additional
638 information is contained within the action estimation vector generated by BoW, and,
639 therefore, there is scope to exploit it to improve the initial annotation. This is exactly
640 what our reasoning system intends to do.



641

642 Figure 4: Blue: Average Cumulative Recognition curve for a number of estimations
 643 from 1 to 13. Red: Recognition rate obtained by our approach depending on the
 644 number of considered action estimates.

645 *b) Performance of the whole framework*

646 The proposed framework integrating AIRS has been tested using the 11 sequences
 647 of WaRo11. Experiments were conducted considering the $N=\{1,3,5,7\}$ most likely
 648 actions estimates – as calculated by CVS - for AIRS analysis. Performance results
 649 are evaluated against the CVS only system in Table 3, where average and
 650 recognition rates per sequence are provided. In addition, they are compared with the
 651 CVS cumulative recognition rate, Figure 4, red.

652 Table 3. Recognition rates obtained using either CVS or the combination of CVS and
 653 AIRS on WaRO11 dataset.

Actor	1	2	3	4	5	6	7	8	9	10	11	Average per action
CVS	35.5%	16.0%	30.0%	58.3%	44.4%	22.2%	40.0%	15.4%	40.0%	16.7%	33.3%	29.4%
CVS+AIRS (n=1)	38.7%	24.0%	30.0%	58.3%	44.4%	22.2%	33.3%	30.8%	60.0%	25.0%	33.3%	35.5%
CVS+AIRS (n=3)	41.9%	28.0%	40.0%	66.7%	44.4%	38.9%	20.0%	30.8%	60.0%	25.0%	33.3%	38.7%
CVS+AIRS (n=5)	64.5%	52.0%	50.0%	75.0%	55.6%	66.7%	40.0%	30.8%	60.0%	25.0%	33.3%	51.9%

CVS+AIRS (n=7)	61.3%	40.0%	60.0%	75.0%	55.6%	66.7%	33.3%	30.8%	40.0%	25.0%	33.3%	51.0%
-------------------	-------	-------	--------------	--------------	--------------	--------------	-------	--------------	-------	--------------	--------------	-------

654

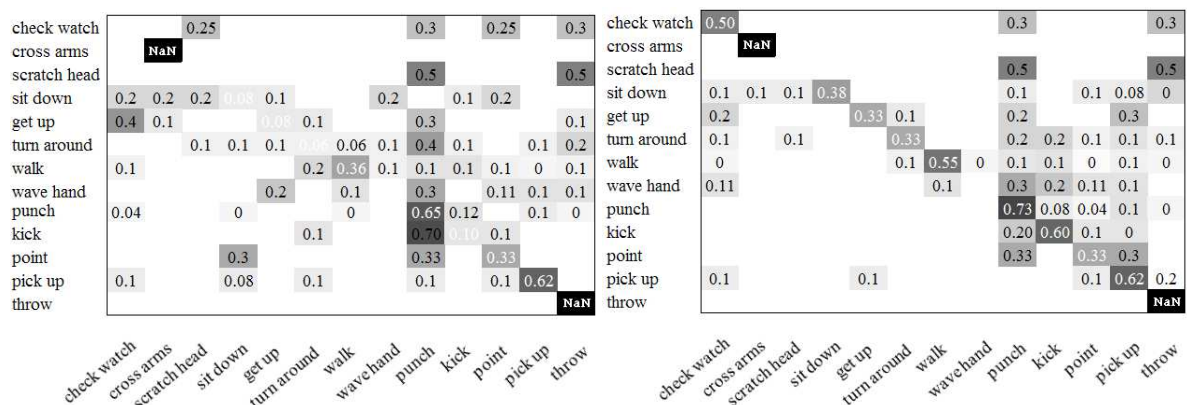
655 These results show a considerable increase of performance due to the inclusion of
656 the reasoning system, i.e. accuracy raises from 29% to 52%, in the best case. Our
657 framework outperforms significantly the CVS system, even for the case where only 1
658 action prediction is considered by the AIRS. Moreover, it can be noticed that
659 accuracy is only rarely deteriorated by reasoning: the system does not seem to
660 suffer from either reasoning errors or over reasoning. Only in sequences 7 and 11
661 performance are either deteriorated or unaffected by the inclusion of the AIRS.
662 Detailed analysis of these two sequences permits to identify, first, absence of
663 continuity or causality between their composing actions and, secondly, a high
664 percentage of unconstrained actions, i.e. actions that are not linked to any other and
665 that can be performed at any instant ('cross arms', 'check watch', 'scratch head').
666 These two factors explain why no effective reasoning can be performed to improve
667 recognition.

668 A more detailed analysis of the AIRS can be obtained by comparing the performance
669 of our approach when varying the number of predictions considered in the action
670 estimate vector. When only considering the most likely action estimate (N=1), the
671 reasoning system is already able to improve on the CVS. This demonstrates the
672 value of one of the AIRS reasoning mechanisms, i.e. activity recognition based on
673 expectations. In this context, the AIRS is comparable to the state-of-art techniques in
674 video-based systems based on simple ontologies and rules.

675 When several action estimates are available, the AIRS's second mechanism, i.e.
676 common sense action validation and the coherent assignation of actions to stories,
677 can be exploited, which leads to deeper reasoning. Performance of the total system

678 – i.e. 38% and 52% for N=3 and 5 estimates, respectively - compared with those
 679 displayed by the ACR – 40% and 57%- shows that the complete reasoning system is
 680 quite efficient at selecting an action among the N best estimates (see Figure 4, red).
 681 Finally, when more estimates are considered, it seems that the added noise prevents
 682 the reasoning system to further improve accuracy, i.e. 51% for N=7.

683 Figure 5 provides confusion matrices with (CVS+AIRS for the best case, i.e. N=5)
 684 and without reasoning (CVS only) to visualise improvement on the recognition rate
 685 per action. For many actions, such as ‘sitting down’, ‘getting up’, ‘turn around’, ‘check
 686 watch’ or ‘kick’, the system is able to move from a recognition rate of almost 0% to a
 687 situation where the action is recognised correctly in a majority of instances. This is
 688 particularly remarkable in the case of ‘sitting down’ where the CVS was trained using
 689 sequences of individuals sitting on the floor, whereas in WaRO11, they sit on a chair.
 690 Such achievement could not have been reached without usage of world and
 691 contextual information. As discussed earlier, recognition rate of an unconstrained
 692 action such as ‘scratch head’ does not benefit from reasoning.



693
 694 Figure 5. Confusion matrices obtained with CVS (left) and CVS+AIRS (right)

695 Table 4: Outputs of CVS (N=5) and AIRS for the first 10 actions of WaRo11 seq. 1



Frames	220-271	271-310	310-344	344-373	373-394
Ground truth	Walk	Pick up	Turn around	Sit down	Get up
CVS 1	Walk	Pick up	Kick	Sit down	Check watch
CVS 2	Kick	Point	Point	Throw	Throw
CVS 3	Point	Throw	Turn around	Check watch	Kick
CVS 4	Wave hand	Scratch head	Pick up	Pick up	Point
CVS 5	Sit down	Sit down	Cross arms	Cross arms	Pick up
AIRS main story	Walk	Pick up	Turn around	Sit down	Get up
					
Frames	394-432	432-1243	1243-1276	1276-1326	1326-1533
Ground truth	Pick up	Sit down	Get up	Pick up	Punch
CVS 1	Pick up	Cross arms	Punch	Pick up	Punch
CVS 2	Get up	Point	Point	Throw	Kick
CVS 3	Throw	Check watch	Kick	Get up	Throw
CVS 4	Scratch head	Scratch head	Pick up	Point	Point
CVS 5	Turn around	Sit down	Throw	Check watch	Check watch
AIRS main story	Turn around	Sit down	Get up	Pick up	Punch

696 Table 4 illustrates the importance of reasoning to improve performance by showing
697 outputs of CVS (N=5) and AIRS for the first 10 actions of sequence 1. When CVS
698 failed to identify the correct actions as its first estimate, AIRS was able to choose the
699 correct annotations among the other 4 estimates, i.e. ‘turn around’ and ‘sit down’
700 actions. Moreover, when none of the CVS outputs was suitable, AIRS managed to
701 correct those estimates by inferring a new action consistent with common sense
702 reasoning – ‘get up’ actions. An error of reasoning occurred in the 6th action, where
703 the AIRS contradicted the correct CVS estimation. This error is explained by the
704 unexpected presence of a second object on the floor, i.e. a pen, which was not

705 known by the DSK. Consequently, the rule imposing that a second object could be
706 picked only after releasing the first one proved invalid.

707 **6. Conclusions**

708

709 We present a novel approach for action recognition based on the combination of
710 statistical and knowledge based reasoning. The inclusion of artificial intelligence
711 strategies, based on common sense, allows outperforming significantly the state of
712 the art technique in computer vision when dealing with realistic datasets. Our main
713 contributions are the creation of the first integrated framework combining computer-
714 vision-based and artificial-intelligence-based action recognition techniques which is
715 fully context and scenario independent, and the implementation of a common sense
716 reasoning schema which outperforms machine learning methodologies.

717 Results are highly encouraging and confirm the validity of our hypothesis: the
718 computer vision community should not focus exclusively on classical statistical
719 reasoning, but should integrate ideas and methodologies from artificial intelligence in
720 order to overcome the limitations of current applications under real-life conditions.

721 **Acknowledgement**

722 This research has been partly supported by the Spanish Ministry of Economy and
723 Competitiveness under the project DREAMS TEC2011-28666-C04-03.

724 **References**

725 Ahmad, M. and Lee, S.-W., 2008. Human action recognition using shape and clg-
726 motion flow from multi-view image sequences. *Pattern Recognition*, 41(7): pp. 2237–
727 2252.

728 Akdemir, U., Turaga, P., Chellappa, R., 2008. An ontology based approach for
729 activity recognition from video. *Proceeding of the 16th ACM international conference*
730 *on Multimedia*, pp.709-712.

731 Ali, A., Aggarwal, J. K., 2001. Segmentation and Recognition of Continuous Human
732 Activity. *IEEE Workshop on Detection and Recognition of Events in Video*.

- 733 Black, M., Yacoob, Y., Jepson, A., Fleet, D., 1997. Learning parameterized models
734 of image motion. IEEE Conf. on Comput. Vis. and Patt. Recog.
- 735 Blackburn, J., Ribeiro, E., 2007. Human motion recognition using isomap and
736 dynamic time warping. Lecture Notes in Computer Science, 4814: pp. 285–298.
- 737 Blank, M., Gorelick, L., Shechtman, E., Irani, M., Basri, R., 2005. Actions as space-
738 time shapes. ICCV.
- 739 Chen, D., Yang, J., Wactlar, H.D., 2004. Towards automatic analysis of social
740 interaction patterns in a nursing home environment from video. Proc. 6th ACM
741 SIGMM Int. Workshop Multimedia Inf. Retrieval, pp. 283–290.
- 742 Chen, W., Fahlman, S.E., 2008. Modeling Mental Contexts and Their Interactions",
743 AAI 2008 Fall Symposium on Biologically Inspired Cognitive Architectures.
- 744 Chen, L. Nugent, C.D., 2009. Ontology-based activity recognition in intelligent
745 pervasive environments. IJWIS 5(4), pp. 410-430.
- 746 Csurka, G., Bray, C., Dance, C., Fan, L., 2004. Visual categorization with bags of
747 keypoints. Workshop on Statistical Learning in Computer Vision, pp. 1–22.
- 748 Dean, M., Schreiber, G., (ed) van Harmelen, F., Hendler, J., Horrocks, I.,
749 McGuinness, D., Patel-Schneider, P., Stein, L., 2011a. OWL Web Ontology
750 Language Reference <http://www.w3.org/TR/2003/WD-owl-ref-20030331/> (last
751 accessed March 2011).
- 752 Dean, M., Schreiber, G. (eds), Bechhofer, S., van Harmelen, F., Hendler, J.,
753 Horrocks, I., McGuinness, D.L. Patel-Schneider, P.F., Stein, L.A., Olin, F.W., 2011b.
754 OWL Web Ontology Language <http://www.w3.org/TR/owl-ref/> (last accessed March
755 2011).
- 756 Duong, T.V, Bui, H.H., Phung, D.Q, Venkatesh, S., 2005. Activity recognition and
757 abnormality detection with the switching hidden semi-markov model. CVPR, pp. 838-
758 845.
- 759 Eagle, N. Singh, P., Pentland, A., 2003. Common sense conversations:
760 understanding casual conversation using a common sense database. Proceedings
761 of the Artificial Intelligence, Information Access, and Mobile Computing Workshop.
- 762 Fahlman, S.E., 2006. Marker-Passing Inference in the Scone Knowledge-Base
763 System. First International Conference on Knowledge Science, Engineering and
764 Management (KSEM'06).
- 765 Fang, C., Chen, J., Tseng, C., Lien, J., 2009. Human action recognition using spatio-
766 temporal classification. Proceedings of the 9th Asian Conference on Computer
767 Vision, pp. 98–109.
- 768 Francois, A.R.J., Nevatia, R., Hobbs, J., Bolles, R.C., 2005. VERL: An Ontology
769 Framework for Representing and Annotating Video Events. IEEE MultiMedia 12(4):
770 pp.76-86.
- 771 Fellbaum, C., 1998. WordNet: An Electronic Lexical Database. MIT Press
- 772 Georis, B., Maziere, M., Bremond, F., Thonnat, M., 2004. A video interpretation
773 platform applied to bank agency monitoring. Proc. 2nd Workshop Intell. Distributed
774 Surveillance System, pp.46–50.

775 Hakeem, A., Shah, M., 2004. Ontology and Taxonomy Collaborated Framework for
776 Meeting Classification. Proc. Int. Conf. Pattern Recognition, pp.219–222.

777 Hobbs, J., Nevatia, R., Bolles, B., 2004. An Ontology for Video Event
778 Representation. IEEE Workshop on Event Detection and Recognition.

779 Ivano, Y., Bobick, A., 2000. Recognition of Visual Activities and Interactions by
780 Stochastic Parsing. IEEE Trans Pattern Analysis and Machine Intelligence .22(8):
781 pp.852–872.

782 Jia, K., Yeung, D., 2008. Human action recognition using local spatio-temporal
783 discriminant embedding. International Conference on Computer Vision and Pattern
784 Recognition, pp. 1–8.

785 Joachims, T., 1998. Text categorization with support vector machines: Learning with
786 many relevant features. ECML.

787 Junejo, I.N., Dexter, E., Laptev, I., Pérez, P., 2008. Cross-view action recognition
788 from temporal self-similarities. ECCV 2008, Part II. LNCS, vol. 5303, pp. 293–306.

789 Kaaniche, M.B., Bremond, F., 2010. Gesture Recognition by Learning Local Motion
790 Signatures. CVPR.

791 Kellokumpu, V., Zhao, G., Pietikäinen, M., 2008. Human activity recognition using a
792 dynamic texture based method. Proceedings of the 19th British Machine Vision
793 Conference, pp. 885–894.

794 Kuipers, B., 1994. Qualitative Reasoning: Modelling and Simulation with Incomplete
795 Knowledge. Cambridge, Mass.: MIT Press.

796 Kovashka, A., Grauman, K., 2010. Learning a hierarchy of discriminative space-time
797 neighborhood features for human action recognition. Proceedings of the International
798 Conference on Computer Vision and Pattern Recognition, pp. 2046–2053.

799 Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T. 2011. HMDB: A Large
800 Video Database for Human Motion Recognition. ICCV.

801 Laptev, I., 2005. On Space-Time Interest Points. International Journal of Computer
802 Vision. 64(2/3): pp. 107–123.

803 Laptev, I., Perez, P., 2007. Retrieving Actions in Movies. ICCV.

804 Laptev, I., Marszalek, M., Schmid, C., Rozenfeld, B., 2008. Learning realistic human
805 actions from movies. Proceedings of the International Conference on Computer
806 Vision and Pattern Recognition, pp. 1–8.

807 Lenat, D., Guha, R.V., 1989. Building Large Knowledge-Based Systems:
808 Representation and Inference in the Cyc project. Addison-Wesley Longman
809 Publishing Co., Inc.

810 Lenat, D., Guha, R.V., Pittman, K., Pratt, D., Shepherd, M., 1990. Cyc: Toward
811 programs with common sense. Commun, ACM, 33(8): pp.30-49.

812 Lewandowski, J., Makris, D., Nebel, J.C., 2010. View and style-independent action
813 manifolds for human activity recognition. Proc. ECCV 6316.

814 Lewandowski, J., Makris, D., Nebel, J.C., 2011. Probabilistic Feature Extraction from
815 Time Series using Spatio-Temporal Constraints. Pacific-Asia Conference on
816 Knowledge Discovery and Data Mining.

817 Liu, J., Ali, S., Shah, M., 2008. Recognizing human actions using multiple features.
818 Proceedings of the International Conference on Computer Vision and Pattern
819 Recognition.

820 Liu, J., Shah, M., 2008b. Learning human actions via information maximization.
821 Proceedings of the International Conference on Computer Vision and Pattern
822 Recognition.

823 McCarthy, J., 1968. Programs with Common Sense. *Semantic Information*
824 *Processing*, Vol. 1, pp. 403–418.

825 McCarthy, J., 1979. Ascribing mental qualities to machines. *Philosophical*
826 *Perspectives in Artificial Intelligence*, pp. 167-195.

827 Makris, D., Ellis, T., Black, J., 2008 *Intelligent Visual Surveillance: Towards Cognitive*
828 *Vision Systems*. *The Open Cybernetics and Systemics Journal*, 2, pp. 219-229.

829 Martinez F., Orrite, C., Herrero, E., Ragheb, H., Velastin, S., 2009. Recognizing
830 human actions using silhouette-based HMM. *Proceedings of the 6th International*
831 *Conference on Advanced Video and Signal Based Surveillance*, pp 43–48.

832 Matikainen, P., Hebert, M., Sukthankar, R., 2010. Representing pairwise spatial and
833 temporal relations for action recognition. *Proceedings of the 11th European*
834 *Conference on Computer Vision*.

835 Minsky M., 1986. *The society of mind*. Simon & Schuster, Inc.

836 Moore, D.J., Essa, I.A., Hayes, M.H., 1999. Exploiting human actions and object
837 context for recognition tasks. *ICCV*, pp 80-86.

838 Natarajan, P., Nevatia, R., 2008. View and scale invariant action recognition using
839 multiview shape-flow models. *Proceedings of the International Conference on*
840 *Computer Vision and Pattern Recognition*, pp. 1–8.

841 Nebel, J.C., Lewandowski, M., Thevenon, J., Martinez, F., Velastin, S., 2011. Are
842 Current Monocular Computer Vision Systems for Human Action Recognition Suitable
843 for Visual Surveillance Applications? *International Symposium on Visual Computing*.

844 Orrite, C., Martinez, F., Herrero, E., Ragheb, H., Velastin, S.A., 2008. Independent
845 viewpoint silhouette-based human action modeling and recognition. *MLVMA*.

846 Philipose, M., Fishkin, K.P., Perkowitz, M., Patterson, D.J., Kautz, H., Hahnel, D.,
847 2004. Inferring activities from interactions with objects. *IEEE Pervasive Computing*
848 *Magazine*, 3(4): pp. 50-57.

849 Richard, S., Kyle, P., 2009. Viewpoint manifolds for action recognition. *EURASIP*
850 *Journal on Image and Video Processing*.

851 Rui, Y. Anandan, P., 2002. Segmenting visual actions based on spatiotemporal
852 motion patterns. *CVPR*.

853 Ryoo, M.S., Aggarwal, J.K., 2009. Spatio-Temporal Relationship Match: Video
854 Structure Comparison for Recognition of Complex Human Activities. *ICCV*.

855 Schuldts, C., Laptev, I., Caputo., B., 2004. Recognizing human actions: A local SVM
856 approach. *ICPR*.

857 Shi, Q., Wang, L. Cheng, L., Smola, A., 2011. Discriminative Human Action
858 Segmentation and Recognition using Semi-Markov Model, *International Journal of*
859 *Computer Vision*, 93(1): pp. 22-32.

860 Shimosaka, M., Mori, T., Sato, T., 2007. Robust Action Recognition and
861 Segmentation with Multi-Task Conditional Random Fields. IEEE International
862 Conference on Robotics and Automation, pp. 3780 - 3786.

863 Ta, A., Wolf, C., Lavoué, G., Baskurt, A., Jolion, J.-M., 2010 Pairwise features for
864 human action recognition. Proceedings of the 20th International Conference on
865 Pattern Recognition.

866 Tapia, E.M, Intille, S., Larson, K., 2004. Activity recognition in the home using
867 simple and ubiquitous sensors. Pervasive, pp. 158-175.

868 Turaga, P., Veeraraghavan, A., Chellappa, R., 2008. Statistical analysis on stiefel
869 and grassmann manifolds with applications in computer vision. International
870 Conference on Computer Vision and Pattern Recognition, pp. 1–8.

871 Vezzani, R., Baltieri, D., and Cucchiara, R., 2010. HMM based action recognition
872 with projection histogram features. Proceedings of the 20th International Conference
873 on Pattern Recognition: Contest on Semantic Description of Human Activities.

874 Vu, V.T., Bremond F., Thonnat, M. 2002. Temporal Constraints for Video
875 Interpretation. 15th European Conference on Artificial Intelligence.

876 Waltisberg, D., Yao, A., Gall, J., Van Gool, L., 2010. Variations of a hough-voting
877 action recognition system. ICPR 2010. LNCS, vol. 6388, pp. 306–312.

878 Wang, L., Suter, D., 2007. Recognizing human activities from silhouettes: Motion
879 subspace and factorial discriminative graphical model. Proceedings of the
880 International Conference on Computer Vision and Pattern Recognition, pp. 1–8.

881 Wang, L., Suter, D., 2007b. Learning and matching of dynamic shape manifolds for
882 human action recognition. IEEE Transactions on Image Processing, 16(6): pp. 1646–
883 1661.

884 Wang, S., Pentney, W., Popescu, A.M., Choudhury, T., Philipose, M., 2007c.
885 Common Sense Based Joint Training of Human Activity Recognizers. Proc.
886 International Joint Conference on Artificial Intelligence.

887 Wang, L. and Suter, D., 2008. Visual learning and recognition of sequential data
888 manifolds with applications to human movement analysis. Computer Vision and
889 Image Understanding, 110(2): pp. 153–172.

890 Santofimia, M.J., Martinez-del-Rincon, J., Nebel, J.C., 2012. WaRo11 Dataset
891 (under development)

892 Weinland, D., Boyer, E., and Ronfard, R., 2007. Action recognition from arbitrary
893 views using 3d exemplars. Proceedings of the 11th International Conference on
894 Computer Vision, 5(7):8.

895 Weinland, D., Ronfard, R., Boyer, E., 2006. Free viewpoint action recognition using
896 motion history volumes. Computer Vision and Image Understanding 104(2-3), pp.
897 249–257.

898 Weinland, D., Özuysal, M., Fua, P., 2010. Making Action Recognition Robust to
899 Occlusions and Viewpoint Changes. ECCV.

900 Yan, P., Khan, S., Shah, M., 2008. Learning 4D action feature models for arbitrary
901 view action recognition. CVPR.

902 Zhang, J., Gong, S., 2010. Action categorization with modified hidden conditional
903 random field. Pattern Recognition, 43(1): pp. 197–203.