

Bayesian Filter based Behavior Recognition in Workflows allowing for User Feedback

Dimitrios I. Kosmopoulos*, Nikolaos D. Doulamis**, Athanasios S. Voulodimos**

Abstract

In this paper, we propose a novel online framework for behavior understanding, in visual workflows, capable of achieving high recognition rates in real-time. To effect online recognition, we propose a methodology that employs a Bayesian filter supported by hidden Markov models. We also introduce a novel re-adjustment framework of behavior recognition and classification by incorporating the user's feedback into the learning process through two proposed schemes: a plain non-linear one and a more sophisticated recursive one. The proposed approach aims at dynamically correcting erroneous classification results to enhance the behavior modeling and therefore the overall classification rates. The performance is thoroughly evaluated under real-life complex visual behavior understanding scenarios in an industrial plant. The obtained results are compared and discussed.

Keywords:

Hidden Markov Models, behavior recognition, workflow, user interaction

1. Introduction

2 The great usefulness of human behavior recognition and understanding in
3 a wide range of applications has attracted the interest of many researchers in

*University of Texas at Arlington, Computer Science and Engineering, Texas, USA

**National Technical University of Athens, School of Electrical & Computer Engineering, Zografou, Greece

Email addresses: dkosmo@ieee.org (Dimitrios I. Kosmopoulos),
ndoulam@cs.ntua.gr (Nikolaos D. Doulamis), thanosv@mail.ntua.gr (Athanasios S. Voulodimos)

4 the areas of computer vision and machine learning. Virtual reality, human-
5 computer interaction, smart environments building and smart monitoring
6 are just a few applications, to which the significance of behavior recognition
7 is indubitable. Especially when it comes to smart monitoring of large-scale
8 enterprises/factories, such as industrial assembly lines, the importance of
9 behavior recognition relates to the safety and security of the staff, to the
10 reduction of costs, to production scheduling, as well as to the quality of
11 the production process. The latter is guaranteed by enforcing adherence to
12 strictly predefined procedures and activities for production or service provi-
13 sion.

14 In most current approaches the goal is either to detect activities, which
15 may deviate from the norm, or to classify some isolated activities. Attempts
16 to address the problem under discussion are encumbered by a number of
17 important hindering factors; the high diversity of the actions and types of
18 behaviors to be recognized is probably the most important one. The complex-
19 ity of static object detection and moving object tracking, with the occlusions
20 and illumination changes, naturally affect adversely approaches that follow
21 the bottom-up paradigm.

22 Despite the above impediments, focusing on monitoring the production
23 line of an industrial plant (such as an automobile manufacturer), which is a
24 fairly structured process, makes modeling of the activities more realistic than
25 in the case of a more unsystematic area of interest, e.g., an airport or a service
26 maintenance system. The former processes are often hierarchically structured
27 as workflows, that comprise sequential tasks. As opposed to isolated action
28 monitoring, the goal here is to monitor activities that occur continuously.

29 This paper proposes innovative workflow recognition schemes for complex
30 industrial environments such as the one of an automobile construction. The
31 main idea is to classify online the visual observations into the available classes
32 (industrial tasks), so that better production monitoring can be achieved. The
33 proposed methodology recognizes the visual tasks as they are captured by the
34 camera overcoming difficulties arising from the complexity of the environment
35 and of the visual process. In addition, we exploit the expert feedback on part
36 of the footage so as to enhance future classification results.

37 The online classification is performed in our case by combining a prob-
38 abilistic theory with supervised time series classifier such as hidden Markov
39 models (HMMs). The sole use of a probabilistic framework cannot solve the
40 problem since we need to know the "a priori" statistics of the visual process.
41 On the other hand, the sole use of HMMs makes the problem only suitable

42 for categorizing pre-segmented sequences, which severely deteriorate the in-
43 dustrial impact, since the expert users are forced to segment a priori the
44 content (see e.g., [1], [2]).

45 Taking these observations into consideration, the work presented in this
46 paper contributes mainly in the following ways:

- 47 • We present a novel online classification framework for distinct behavior
48 recognition in visual workflows. The proposed framework is based on
49 HMMs and Bayesian filtering and exploits prior knowledge.
- 50 • We also propose an approach for improving the supplied results by
51 allowing interaction with the user, who may provide relevance feedback.
52 Two different neural network based schemes are introduced: non-linear
53 and recursive, the latter allowing for significant error decrease using a
54 small number of training samples.

55 Furthermore, in contrast to most mainstream bottom-up approaches we
56 employ features at the image level, bypassing the error-prone object detection
57 and tracking steps.

58 The rest of this paper is organized as follows: After surveying the related
59 literature regarding behavior understanding in section 2, we formulate the
60 problem that we are proposing to solve in section 3. Then describe the pro-
61 posed task representation in section 4 and the proposed online classification
62 framework in section 5. Section 6 focuses on the neural network based rec-
63 tification schemes. The experimental setup and the outcoming results are
64 described and analyzed in section 7. Finally, section 8 concludes the paper
65 with a summary of the findings.

66 2. Related Work

67 Event detection and especially human action recognition has been the
68 focus of interest of computer vision and machine learning communities for
69 years, mostly as isolated activities and not as part of a continuous process.
70 A variety of methods has addressed these problems, including semi-latent
71 topic models [3], spatial-temporal context [4], optical flow and kinematic fea-
72 tures [5], and random trees and Hough transform voting [6]. Wada et al.
73 [7] employ Non-deterministic Finite Automaton as a sequence analyzer to
74 present an approach for multi-object behavior recognition based on behavior
75 driven selective attention. Other works focus on more specific domains, e.g.,

76 event detection in sports [8, 9], retrieving actions in movies [10], human ges-
77 ture recognition (using Dynamic Time Warping [11] and Time Delay Neural
78 Networks [12]), and automatic discovery of activities [13]. Comprehensive
79 literature reviews regarding isolated human action recognition can be found
80 in [14, 15].

81 One of the key functionalities of any machine learning model (classifier)
82 suitable for application in visual behavior understanding is the ability to ex-
83 tract the *signature* of a behavior from the captured visual input. The key
84 requirements when designing such a classifier is (a) to support task execution
85 in various time scales, since a task or parts of it may have variable duration;
86 and (b) to support stochastic processes, because of the task intra-class vari-
87 ability and noise.

88 A very flexible framework for stochastic classification of time series is
89 the HMM (see e.g., [16]). It can be easily extended to handle outliers (see
90 e.g., [17]) and to fuse multiple streams (e.g., [18]). It is very efficient for
91 application in previously segmented sequences (see e.g., [19], [20]), however
92 when the boundaries of the sequence that we aim to classify are not known in
93 advance, the search space of all possible beginning and end points make the
94 search very inefficient [21]. A typical way to treat this problem is given in
95 [22], where a dynamic programming algorithm of cost T^3 , is used to perform
96 segmentation and classify then the segments; however the cost is restrictive
97 in real applications.

98 In the past there have been some efforts to exploit the hierarchical struc-
99 ture of some time series, e.g., by using the hierarchical HMMs [23]. Each
100 state is considered to be a self-contained probabilistic model (an HHMM).
101 Examples of such approaches can be found in [24], where the workflow in a
102 hospital operating room is described. Another approach is the layered hidden
103 Markov model (LHMM) (see [25]), which consists of N levels of HMMs where
104 the HMMs on level $N + 1$ corresponds to observation symbols or probability
105 generators at level N . Every level i of the LHMM consists of K_i HMMs
106 running in parallel. In that work a LHMM is used for event identification in
107 meetings. In [26] structure learning in HMMs is addressed in order to obtain
108 temporal dependencies between high-level events for video segmentation. An
109 HMM models the simultaneous output of event-classifiers to filter the wrong
110 detections.

111 In many workflows, such as in industrial production where a sequence of
112 different tasks has to be completed, the execution of a task means that it
113 will not appear again in the same workflow. Therefore the whole history of

114 tasks must be kept in memory to exclude false positives and the Markovian
115 property is obviously not applicable. Thus, the above approaches have an
116 inherent problem to describe such workflows.

117 The Echo State Network (ESN), see, e.g., [27], could be a promising
118 method for online classification of workflow time series, because it does not
119 make any explicit Markovian assumption. However, it was shown in [28] that
120 it effectively behaves as a Markovian classifier, i.e., recent states have a far
121 larger influence on the predicted state. The ESN has already been used in a
122 work using the same dataset that we are using [29]. However, their results
123 are not directly comparable to ours, since the features they are using are
124 different.

125 In the proposed work we aim to alleviate the problems of online task
126 segmentation and we by-pass the erroneous Markovian assumption by ap-
127 proximating the probability of a label sequence (for each incoming frame)
128 given the whole observation history. To this end we employ Bayesian filter-
129 ing. The related techniques for online behavior recognition have not been
130 adequately investigated in the literature. In the works of [30] and [31] Rao
131 Blackwell particle filters were used along with a dynamic Bayesian network
132 for tracking of hierarchical events. In our work we make no assumptions
133 about event hierarchy and the representation that we adopt is by definition
134 very simple, so resorting to Rao Blackwell filters for state space reduction
135 is not needed. The work in [32] notices the utility of particle filters in com-
136 bination with an HMM, however it seeks to perform observation prediction,
137 which is different from our classification problem.

138 A scheme using an HMM in combination with a particle filter was pre-
139 sented in [33] to model well-log data. More specifically, by using a single
140 (modified) HMM, consisting of hidden and measurable states and with the
141 help of a particle filter the sequence of HMM states was extracted. In con-
142 trast to [33], we propose online classification, i.e., to find which of the tasks
143 has generated the current and all previous observations, using several HMMs
144 (each of them modeling a separate task) and prior knowledge about task
145 execution. Each visual task has to be modeled by a separate HMM due to
146 their dynamic characteristics and due to their high complexity. Our work
147 contributes by showing how to execute online multi-class classification by
148 defining appropriately the proposal function, the model of prior knowledge
149 and the integration of HMM models with a particle filter.

150 Regarding relevance feedback, it is a common approach for automati-
151 cally adjusting the response of a system regarding information taken from

152 user’s interaction [34]. Originally, it has been developed in traditional infor-
153 mation retrieval systems [35], [36], but it has been now extended to other
154 applications, such as surveillance systems [37], [38]. Relevance feedback is
155 actually an online learning strategy which re-weights important parameters
156 of a procedure in order to improve its performance. Re-weighting strategies
157 can be linear or non-linear relying either on heuristic or optimized method-
158 ologies [34]. Linear and heuristic approaches usually adjust the degree of
159 importance of several parameters involved in the selection process. Instead,
160 non-linear methods adjust the applied method itself using function approx-
161 imation strategies [39]. In this direction neural network models have been
162 introduced as non-linear function approximation systems. However, such
163 approaches have been applied for information retrieval systems instead of
164 surveillance applications. It is clear that it is not straightforward to move
165 from one type of application to another due to the quite different require-
166 ments of both applications.

167 A comprehensive review regarding algorithms of relevance feedback in
168 image retrieval has been provided in [40]. In this paper, the authors lay em-
169 phasis on comparing different techniques of relevance feedback with respect
170 to the type of training data, the adopted organization strategies, the simi-
171 larity metrics used, the implemented learning strategies and finally the effect
172 of negative samples in the training performance. However, the compared
173 methods in [40] are designed in a static, non-dynamic framework. Instead, in
174 real life applications, there exists a non-linear and time varying relationship
175 between the input feature vectors and the target output states. Therefore,
176 we need to introduce dynamic approaches for an efficient relevance feedback
177 implementation.

178 In [41] transfer learning is adopted to address the aforementioned diffi-
179 culties. The method is applied to detect events in Web videos [41]. With
180 transfer learning, we can use auxiliary data from known classification prob-
181 lems, different from the user’s target query, to decrease the amount of data
182 needed to be fed back. The drawback of transfer learning is that it is actu-
183 ally assumed that the previous known environmental conditions, over which
184 the classifier has been trained with, are similar with the current character-
185 istics of the environment. This implies stationary environments. Instead,
186 the performance of the relevance feedback in real life non stationary condi-
187 tions is highly deteriorated. This is addressed in this paper by introducing
188 an adaptable mechanism able to dynamically adjust the trained non-linear
189 relationship using few samples that represent the current environmental con-

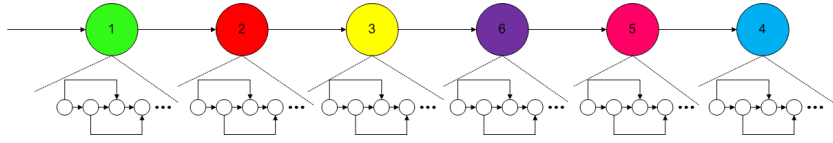


Figure 1: Graph representing the temporal relationships between the tasks of a typical workflow (the order of execution of tasks is an example; there are many possible permutations). Colored nodes represent tasks and white smaller nodes stand for intra-task states, which are the states of the HMM that models the respective task.

190 ditions along with a minimum degradation of the already obtained previous
 191 knowledge of the classifier.

192 3. Problem formulation

193 As stated in section 1, we focus on detecting and recognizing visual tasks
 194 in complex industrial processes (workflows), being executed in an automo-
 195 bile production line. The visual tasks are recognized from visual cues being
 196 captured from a set of cameras.

197 A workflow is a process that happens repetitively and consists of a se-
 198 quence of discrete tasks. The order in which tasks appear matters, however
 199 permutations are allowed in some cases (which have to be learned). Tasks
 200 may have different durations, as a result of the natural differences in work-
 201 ers' productivity. The definition of tasks stems from domain knowledge. An
 202 example of such a task is: "A worker picks part1 from rack1 and places it on
 203 the welding cell".

204 A graph that presents the hierarchy of tasks that compose a workflow
 205 and their internal states is given in Fig. 1. Each workflow is composed of
 206 tasks and each task is modeled by a separate HMM.

207 Our goal is to determine which tasks are executed and when, given in-
 208 stances of workflows, which are described by sequences of visual observations.
 209 Let us denote as \mathbf{o}_t the visual observation vector at the t time instance, or
 210 in discrete domain the frame number t . These visual observations are de-
 211 scriptors (features) extracted by processing the pixel values of frame t . The
 212 visual observations are described in subsection 4.1.

213 Our goal can be alternatively expressed as the classification of each frame
 214 t to one of the L available classes, i.e., different industrial tasks. Let us denote
 215 as $x_t = l_t$ the state vector including the label l_t from the L classes (tasks)
 216 that has to be assigned to frame t . Our goal is to calculate the posterior

217 $p(x_{0:t}|\mathbf{o}_{1:t})$ at every step, given the measurements (visual observations) up to
 218 that step.

219 The notation $x_{0:t}$, which will be used in the following section, accumulates
 220 all vectors x_i with $i = 0, \dots, t$, that is $x_{0:t} = (x_0..x_t)$. Similarly, the notation
 221 $\mathbf{o}_{1:t}$ accumulates all observation descriptors up to time t , that is, $\mathbf{o}_{1:t} =$
 222 $(\mathbf{o}_1..\mathbf{o}_t)$.

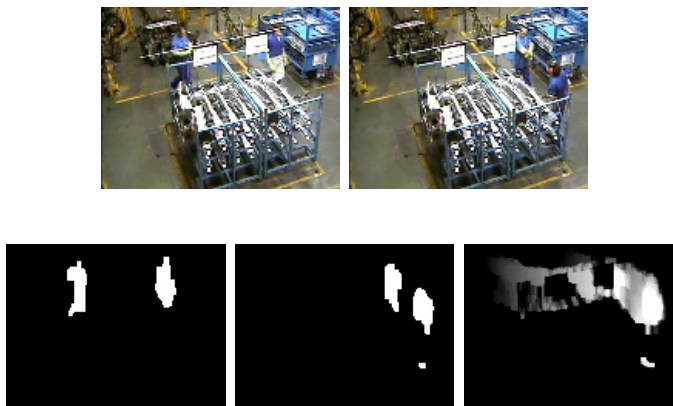


Figure 2: Two keyframes (first row), the respective background subtraction images and the extracted PCH image (second row)

223 4. Task modeling

224 In the following we describe how we represent the tasks. In subsection
 225 4.1 we present means of representation of each frame, while in subsection 4.2
 226 we present the HMM for modeling time series.

227 4.1. Visual Observations

228 One of the key challenges real-time action recognition systems are con-
 229 fronted with concerns selection of appropriate features for representing the
 230 observed raw data. The ideal features should describe different actions ac-
 231 curately, with high discrimination capability, and should be efficiently calcu-
 232 lated. Ideally, these features should also provide a hierarchical representation
 233 scheme (coarse to fine) so that a desirable, application-wise, trade-off between
 234 representation capabilities and computational complexity can be reached.

235 The employment of features directly extracted from the video frames
 236 has the significant advantage of *obviating the need of detecting and tracking*
 237 *the salient scene objects*, a task which is notoriously difficult in cases of
 238 occlusions, target deformations, illumination changes etc. Thus, by using
 239 such an approach, the intermediate levels of semantic complexity, as met in
 240 typical *bottom-up* systems, are completely bypassed. For this purpose, either
 241 local or holistic features (or both [42]) may be used. *Holistic features* such
 242 as Pixel Change History (PCH) images [43] remedy the drawbacks of local
 243 features, while also requiring a much less tedious computational procedure
 244 for their extraction. A very positive attribute of such representations is that
 245 they can easily capture the history of a task that is being executed.

The PCH value of a pixel is defined as:

$$P_{\varsigma,\tau}(x,y,t) = \begin{cases} \min(P_{\varsigma,\tau}(x,y,t-1) + \frac{255}{\varsigma}, 255) \\ \text{if } D(x,y,t) = 1 \\ \max(P_{\varsigma,\tau}(x,y,t-1) - \frac{255}{\tau}, 0) \\ \text{otherwise} \end{cases} \quad (1)$$

246 where $P_{\varsigma,\tau}(x,y,t)$ is the PCH for a pixel at (x,y) , $D(x,y,t)$ is the binary
 247 image indicating the foreground region, ς is an accumulation factor and τ
 248 is a decay factor. By setting appropriate values to ς and τ we are able to
 249 capture pixel-level changes over time (see Fig. 2). These images can then
 250 be transformed to a vector-based representation using moments such as the
 251 Zernike moments (see, e.g., [19]).

252 4.2. The HMM framework and its drawbacks

253 A common approach for stochastically modeling time series is to use hid-
 254 den Markov models (HMMs). A hidden Markov model consists of states,
 255 transitions, observations and probabilistic behavior, and is formally defined
 256 as a tuple $\lambda = \langle \mathbf{Q}, \mathbf{A}, \mathbf{B}, \pi \rangle$ satisfying the following conditions:

- 257 • $\mathbf{Q} = \{q_1, \dots, q_S\}$ is a finite set of S states. In our case, the number
 258 of states is an indication of the order (complexity) of the stochastic
 259 representation.
- 260 • \mathbf{A} is the transition matrix, which represents the transition probabilities
 261 between states.
- 262 • \mathbf{B} is the observation matrix, which represents the observation proba-
 263 bility given the state.

- 264 • π represents the probability of each state at the beginning of the se-
265 quence.

266 A supervised training algorithm is used to obtain the parameters λ of the
267 HMM. The training set is formed using representative samples of industrial
268 tasks which have been manually classified to one of the L available classes.
269 This implies that we need first to annotate the tasks, exploiting, for exam-
270 ple, the experience of industrial engineers. We also need to identify the start
271 and finish times for each industrial workflow even during the testing phase,
272 which is a burden for a real-life exploitation of the algorithm in industrial
273 environments. In real-world scenarios the starting and ending times of tasks
274 are usually unknown. Therefore, HMM modeling can not be used for online
275 recognition of the tasks. This is because online classification requires search-
276 ing in the space of possible beginning and end points to perform Viterbi
277 matches in order to find the optimally fitting sequence [16].

278 Assuming that tasks' appearance follow Markovian behavior (the condi-
279 tional probability distribution of future tasks depends only upon the present
280 task; that is, given the present, the future does not depend on the past) it
281 is possible to perform online classification by applying techniques such as
282 hierarchical (HHMM) and Layer hidden Markov models (LHMM) [23], [25].
283 However, such assumptions are not true in a real-world industrial environ-
284 ment, since the processes considered are structured. Usually, in a real-world
285 production environment, the current execution of a task will affect the exe-
286 cution of future tasks, i.e., a task may be executed only once in a workflow.

287 All the above imply that the use of a conventional HMM for stochasti-
288 cally classifying industrial tasks is very inefficient, especially for real world se-
289 quences, which typically contain several thousands of frames. An exhaustive
290 search for all possible combinations would be therefore practically prohibitive
291 from a computational point of view. Hence, for an online classification frame-
292 work, we need to identify the time boundaries, that is the start and finish
293 times of an industrial task, which are part of a workflow. For this reason,
294 an alternative methodology is required, which constitutes one of the main
295 contributions of this paper.

296 To this end we propose in the following an approximate, though very effi-
297 cient, method, which endows the HMM with online classification capabilities.

298 **5. The Bayesian Filter based Classification Framework**

299 A method for online recognition of industrial tasks in visual workflows
 300 based on Bayesian filters is proposed here. We assume that we are not aware
 301 of the start and finish times of the tasks.

302 As stated in section 3, our goal is to determine which tasks are executed
 303 and when, given instances of workflows. In other words, our goal is to cal-
 304 culate the posterior probability $p(x_{0:t}|\mathbf{o}_{1:t})$ for every frame t . Estimation of
 305 the posterior probability $p(x_{0:t}|\mathbf{o}_{1:t})$ is a much more complex process than
 306 estimating the posterior $p(x_t|\mathbf{o}_t)$, since in the former case, the probability
 307 depends on the classification results of the previous frames.

308 One possible method to calculate $p(x_{0:t}|\mathbf{o}_{1:t})$ is by employing a Bayesian
 309 filter. The solution for the Bayesian filter is commonly expressed as:

$$p(x_{0:t}|\mathbf{o}_{1:t}) = p(x_{0:t-1}|\mathbf{o}_{1:t-1}) \frac{p(\mathbf{o}_t|x_{0:t}, \mathbf{o}_{1:t-1}) p(x_t|x_{0:t-1}, \mathbf{o}_{1:t-1})}{p(\mathbf{o}_t|\mathbf{o}_{1:t-1})} \quad (2)$$

310 Equation (2) is actually a recurrent expression of the probability $p(x_{0:t}|\mathbf{o}_{1:t})$
 311 with respect to the previous estimates $p(x_{0:t-1}|\mathbf{o}_{1:t-1})$ up to time $t - 1$. How-
 312 ever, the main difficulty in calculating $p(x_{0:t}|\mathbf{o}_{1:t})$ using equation (2) is the
 313 fraction term on the right part. To estimate this term, we need additional
 314 knowledge regarding the distribution of visual observations of image frame
 315 $t - 1$ being aware of the class (i.e., task) that this frame belongs to.

316 One possible way to estimate the additional knowledge, required for the
 317 online classification framework, is to exploit a supervised classifier as the
 318 HMM, described in subsection 4.2. For this reason, we combine the HMM
 319 with the probabilistic framework, indicated by equation (2) to achieve online
 320 recognition of industrial tasks, disencumbered from the requirement to know
 321 start and finish times in advance.

322 To estimate the fraction term of the right part of equation (2) we proceed
 323 as follows. First, the term $p(\mathbf{o}_t|\mathbf{o}_{1:t-1})$ is independent of the class to which the
 324 current frame should be assigned to, so it can be omitted from the following
 325 calculations.

326 Second, it is reasonable to assume that the observation \mathbf{o}_t depends only on
 327 the current task x_t , so we simplify $p(\mathbf{o}_t|x_{0:t}, \mathbf{o}_{1:t-1})$ to $p(\mathbf{o}_t|x_t)$. We propose to
 328 calculate this probability by using the observation model of the HMM, which
 329 is learned offline for each HMM state. Third, for the term $p(x_t|x_{0:t-1}, \mathbf{o}_{1:t-1})$
 330 we propose an alternative expression, which is the $p(x_t, c_t|x_{0:t-1}, \mathbf{o}_{1:t-1})$ or

331 more simply $p(x_t, c_t|x_{0:t-1})$; the latter holds because if the task history is
 332 known then the observation history does not affect the appearance of the
 333 next task. The variable c_t is a boolean stochastic variable, which becomes
 334 true if the task label changes from $t - 1$ to t and false if the task remains the
 335 same. $p(x_t, c_t|x_{0:t-1})$ gives the probability that the task in current frame t
 336 has the label x_t and there is (or there is not) a switch to a new task, provided
 337 that the sequence of all previous task labels (task history) is known. More
 338 details about this probability can be found in sub-section 5.1.2.

339 Under these assumptions equation (2) becomes:

$$p(x_{0:t}|\mathbf{o}_{1:t}) \propto p(x_{0:t-1}|\mathbf{o}_{1:t-1}) \cdot p(\mathbf{o}_t|x_t) p(x_t, c_t|x_{0:t-1}) \quad (3)$$

340 As observed in equation (3), the posterior probability $p(x_{0:t}|\mathbf{o}_{1:t})$ is pro-
 341 portional to a) the recurrent term $p(x_{0:t-1}|\mathbf{o}_{1:t-1})$, b) the probability $p(\mathbf{o}_t|x_t)$,
 342 which is estimated through the HMM model that captured the supervised
 343 knowledge of the task execution with regard to the visual observations and
 344 c) the $p(x_t, c_t|x_{0:t-1})$, which expresses our a priori knowledge about task du-
 345 ration and transition from one state to another.

346 It is clear, therefore, that the estimation of the probability $p(x_{0:t}|\mathbf{o}_{1:t})$
 347 requires the a priori knowledge about task duration and transition, as well
 348 as the distribution of the visual observations (e.g., visual descriptors) with
 349 respect to the task that a frame belongs to. However, another difficulty of
 350 solving equation (3) is that it involves dependencies from previous frame
 351 observations (visual descriptors) and classification (frame assignment to one
 352 of the L available classes).

353 To handle the dependencies of the posterior probability $p(x_{0:t}|\mathbf{o}_{1:t})$ with
 354 the previous frame states (e.g., frame classification), we need first to intro-
 355 duce a list of hypotheses and then to validate them under a probabilistic
 356 framework. A common approach for performing that is through the usage of
 357 Particle Filters theory, which is a method for estimating the importance of
 358 a hypothesis according to a set of observed data.

359 5.1. Particle Filter Driven by the Hidden Markov Model

360 Let us assume that we have a set of N available hypotheses (particles).
 361 A hypothesis describes a particular combination of the classes that the pre-
 362 vious frames have been assigned to. For example, a hypothesis is that the
 363 first frame belongs to the second task, the second frame to the same task,
 364 while the third to the first task, etc. Every hypothesis is evaluated through

365 the Bayesian filters, which are estimation methods based on simulation and
 366 previous observations [44], [45].

367 Weights are associated to the hypotheses, expressing the significance de-
 368 gree to the modeling process. Therefore given N hypotheses we have N
 369 weighted particle trajectories $\left\{x_{0:t-1}^{(n)}, w_{0:t-1}^{(n)}\right\}_{n=1}^N$. Each of these trajectories
 370 approximates the posterior probability $p(x_{0:t-1}|\mathbf{o}_{1:t-1})$ up to time $t - 1$.

371 Let us assume that the particle trajectories up to time $t - 1$ are known.
 372 Then, we can compute the N particles $\left\{x_t^{(n)}\right\}_{n=1}^N$ which are combined with
 373 the previous trajectories to form $\left\{x_{0:t}, w_{0:t}^{(n)}\right\}_{n=1}^N$ to approximate the pos-
 374 terior $p(x_{0:t}|\mathbf{o}_{1:t})$ up to time t . In particular, the current weight $w_t^{(n)}$ for
 375 the n^{th} hypothesis at the current frame t is estimated through the following
 376 distribution:

$$w_t^{(n)} = p(\mathbf{o}_t|x_t^{(n)}) \quad (4)$$

377 Equation (4) means that we can estimate the weights $w_t^{(n)}$ if we know a
 378 hypothesis, i.e., we know the class x_t to which frame t belongs. The pdf in
 379 equation (4) derives from the supervised learning of the HMM. The hypoth-
 380 esis about the value of x_t requires a priori knowledge regarding task duration
 381 and order of tasks. The hypothesis is generated by sampling the distribution
 382 $p(x_t, c_t|x_{0:t-1})$ which is learned offline (see subsection 5.1.1).

383 5.1.1. Estimation of the Observation Probability

384 The observation probability $p(\mathbf{o}_t|x_t)$ depends not only on the currently
 385 executed task but also on the state q of the associated HMM, so it should
 386 be fully written as: $p(\mathbf{o}_t|x_t, q_t)$. Here the HMM state that maximizes the
 387 observation probability is selected for each task.

388 At this point it should be noted that the hidden system state space (cur-
 389 rently executed task) is one-dimensional and discrete, with low number of
 390 possible states (equals the number of possible tasks). Therefore the method
 391 is very efficient and a relatively low number of particles is required.

392 5.1.2. A Priori Knowledge

393 Our knowledge about the task order as well as the task duration can be
 394 used to estimate the distribution $p(x_t, c_t|x_{0:t-1})$, where we recall that c_t is

395 a boolean stochastic variable, which becomes true if the task label changes
 396 from $t - 1$ to t and false if the task remains the same.

397 Using the Bayes rule we can conclude to:

$$p(x_t, c_t | x_{0:t-1}) = p(x_t | x_{0:t-1}, c_t) \cdot p(c_t | x_{0:t-1}) \quad (5)$$

398 The term $p(c_t | x_{0:t-1})$ depends only on the duration d of the current task.
 399 In other words $p(c_t | x_{0:t-1}) \equiv p(c_t | x_{t-1}, d)$. A common approach for model-
 400 ing $p(c_t | x_{t-1}, d)$ is to use a Gaussian mixture model of the respective task
 401 duration, which can be learned offline. Thus, we have:

$$p(c_t | x_{t-1}, d) = \sum_{i=1}^K m_i N(\mu_i, \sigma_i) \quad (6)$$

402 where K the number of mixture components and m_i, μ_i, σ_i the prior, mean
 403 and standard deviation of the i^{th} component.

The other term of equation (5) is given by

$$p(x_t | x_{0:t-1}, c_t) = \begin{cases} 0 & \text{if } c_t = \text{true} \\ 1 & \text{if } c_t = \text{false} \end{cases} \quad (7)$$

404 in the case that $x_t = x_{t-1}$,

405 and by:

$$p(x_t | x_{0:t-1}, c_t) = \begin{cases} 0 & \text{if } c_t = \text{false} \\ T_{x_{0:t}} & \text{if } c_t = \text{true} \end{cases} \quad (8)$$

406 in the case that $x_t \neq x_{t-1}$. $T_{x_{0:t}}$ is the probability of a path in a decision tree
 407 describing the possible transition paths between tasks. The tree is described
 408 in the following.

409 Assuming that the task with value m is possible to appear in the i -th
 410 order, we may denote $x(i) = m$. There will be a node in the i -th level of
 411 the tree with value equal to m . Given that there are several tasks that may
 412 follow that task directly after its completion, i.e., $x(i+1)$ may take n values,
 413 in the tree the node with value $x(i) = m$ will have n descendants, with these
 414 values.

415 The root of the tree is defined to be a virtual node (with no associated
 416 task value), while its children indicate the tasks that may start the workflow.
 417 Additionally, each link that connects a parent node with value $x(i)$ to a child
 418 node with value $x(i+1)$ has an associated value $l(x(i+1), (x(i)))$, which

419 indicates the probability of occurrence of the task $x(i+1)$ directly after $x(i)$
 420 is finished. A complete workflow is represented by a path connecting the root
 421 with any of the tree leaves. Given a path P , the probability of $p(P)$ is given
 422 by $P(p) = \prod_{pathlinks} l(x(i+1), x(i))$, which is the product of the associated
 423 probabilities for all the links in the path.

424 Such a tree can be learned by using a training set of full workflows,
 425 and therefore the "legitimate" paths and their probabilities can be specified.
 426 More specifically, for each parent node we find the possible successors (de-
 427 scendants) and based on the normalized frequency that a specific child is
 428 selected as next task, we assign a probability value to the connecting link.
 429 Since the history of all previous tasks is maintained by using such a tree, we
 430 do not rely on the Markovian assumption.

431 Algorithm 1 presents the steps of the proposed online learning classifi-
 432 cation framework that combines stochastic modeling and HMMs. The pro-
 433 posed framework is depicted in Fig. 3 as dynamic Bayesian network, where
 434 all dependencies are graphically presented.

435 6. A Neural Network-Based Rectification Scheme

436 The main drawback of the aforementioned probabilistic approach is that
 437 the observation probability $p(\mathbf{o}_t|x_t)$ may sometimes give rise to the wrong
 438 task as a result of the EM-based learning, which can be trapped in local
 439 maxima. This in turn may result in some particles taking a wrong "trajec-
 440 tory". To address this difficulty, we present in the following a rectification
 441 framework able to automatically adjust the $p(\mathbf{o}_t|x_t)$ (the link {3} in the
 442 graph of Fig. 3) according to user's feedback. The rectification strategy
 443 is based on the usage of a dynamic non-linear classifier, which is able to
 444 learn the current user's feedback, as expressed by a small set of selected
 445 relevant/irrelevant data, while simultaneously provide a minimum degrada-
 446 tion of the previous obtained knowledge. Since the rectification function is
 447 expected to be non-linear we based our approach on neural networks.

448 At this point we should mention that any method capable of handling non-
 449 linear functional relationships could probably substitute our neural network
 450 based approach. Alternatives that could have been addressed include non-
 451 linear (kernel) regression [46] or random forests [47].

Algorithm 1 Proposed Method

```
{OFFLINE TRAINING}
{Decision tree learning}
DecisionTree = Learn(AllTaskPaths)
{Supervised task learning through HMM}
for  $s = 1$  to NumberOfTasks do
   $\langle \mathbf{Q}_s, \mathbf{A}_s, \mathbf{B}_s, \pi_s \rangle = \text{TrainHMM}(\text{AllTaskTimeSeries})$ 
end for
{ONLINE CLASSIFICATION FRAMEWORK}
while ( $F = \text{AcquireFrame}()$ )  $\neq$  NULL do
   $\mathbf{o}_t = \text{ProcessFrame}(F)$ ; {extraction of visual observations (features) by
  processing the current captured frame}
  {for every Hypothesis do}
  for  $n = 1$  to  $N$  do
     $x_t^{(n)} = \text{Sample } p(x_t, c_t | x_{0:t-1}^{(n)})$ 
    {get HMM state  $q_t$  that maximizes observation probability}
     $q_t = \arg \max_q \{p(o_t | x_t, q)\}$ 
    Weight  $\mathbf{x}_t^{(n)}$  by the following:
    
$$w_t'^{(n)} = p(\mathbf{o}_t | x_t^{(n)}, q_t) \tag{9}$$

  end for
  for  $n = 1$  to  $N$  do
    Normalize the weights by:
    
$$w_t^{(n)} = \frac{w_t'^{(n)}}{\sum_{n_1=1}^N w_t'^{(n_1)}} \tag{10}$$

  end for
  Switching-state particles with low weight are set back to previous state.
  for  $n = 1$  to  $N$  do
    Update  $p(x_{0:t}^{(n)} | o_{1:t})$  {use the recurrent framework described in eq. 3 }
  end for
  The winner is the particle  $n_0 : p(x_{0:t}^{(n_0)} | o_{1:t}) \geq p(x_{0:t}^{(n)} | o_{1:t}), \forall n \in \{1, \dots, N\}$ 
end while
```

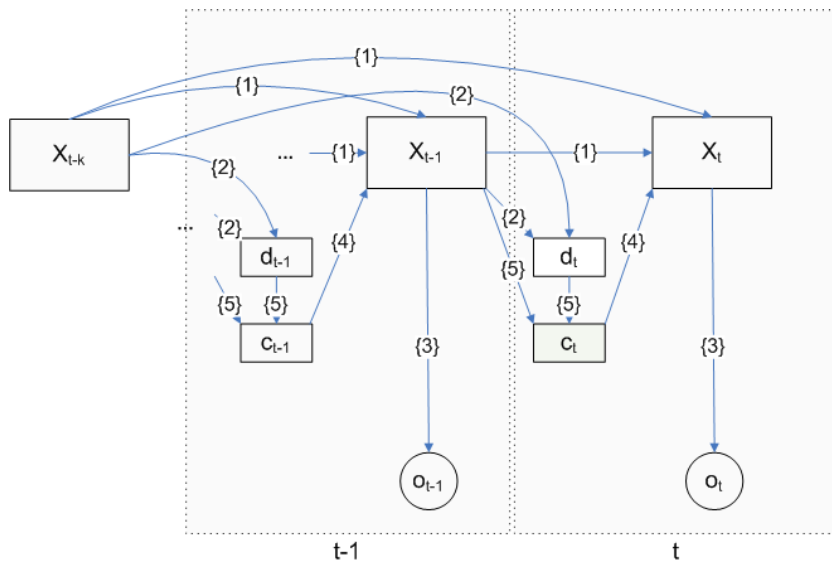


Figure 3: The proposed scheme represented as a dynamic Bayesian network, where rectangles correspond to discrete values and circles to continuous values. The digit-annotated edges represent the dependencies in our framework as follows: {1}: the dependency of the current state x_t on the state history {2}: state duration given the previous states {3}: $p(o_t|x_t)$, derived from HMM. {4} and {1} represent the equations (7,8). {5} represents the equation (6).

452 6.1. The Non-linear Model

453 Let us denote as S a set that contains the selected samples by the user.
 454 The set $S = \{\dots(\mathbf{p}_i, \mathbf{d}_i)\dots\}$ contains pairs of the form $(\mathbf{p}_i, \mathbf{d}_i)$, where as
 455 \mathbf{p}_i we indicate the observation probability vector, generated by the HMM,
 456 the elements of which express the probability of the corresponding frame to
 457 belong to one of the, say, M available classes. Vector \mathbf{d}_i indicates the ideal
 458 probabilities for the i^{th} sample. Variable \mathbf{d}_i is an indicator vector meaning
 459 that all its elements will be zero apart from one which is equal to one. This
 460 element indicates the class that this task belongs to. Assuming the existence
 461 of a non-linear function able to correct the erroneous classifications of the
 462 HMM, we can derive the following equation

$$\mathbf{d}_i = \underline{f}(\mathbf{p}_i) \quad (11)$$

463 where $\underline{f}(\cdot)$ is a vector function indicating the non-linear relationship between
 464 \mathbf{p}_i and \mathbf{d}_i .

465 The main difficulty with the previous equation is the fact that vector
 466 function $f(\cdot)$ is actually unknown. Additionally, the non-linear relationship
 467 dynamically changes under different conditions and camera system modifica-
 468 tion. To address the aforementioned difficulties, we introduce a feed forward
 469 neural network model able to accurately approximate the unknown vector
 470 function $\underline{f}(\cdot)$ with a certain degree of accuracy. In this case, equation (11) is
 471 written as follows:

$$\mathbf{d}_i = \underline{f}_{\mathbf{w}}(\mathbf{p}_i) \quad (12)$$

472 The main difference between equations (11) and (12) is the introduction
 473 of the vector weight \mathbf{w} . This means that different parameters (weights) of
 474 the network yields different performance of the adaptable classifier. Vector \mathbf{w}
 475 includes all the parameters (weights) of the non-linear neural network-based
 476 classifier.

477 To estimate the weights \mathbf{w} we need to apply a training algorithm, which
 478 actually minimizes the mean square error among all selected from the expert
 479 user data (task sequences) and the respective output of the network when a
 480 particular set of weights is applied. That is,

$$\mathbf{w} = \arg \min_{\text{forall } \mathbf{w}} \epsilon = \arg \min_{\text{forall } \mathbf{w}} \sum_i (\underline{f}_{\mathbf{w}}(\mathbf{p}_i) - \mathbf{d}_i)^2 \quad (13)$$

481 The back propagation algorithm [48] can provide a solution to this non-
 482 linear minimization problem. In our experiments, we select a small neural
 483 network structure of few hidden neurons and one hidden layer. In this case,
 484 we try to minimize the number of neural networks parameters, that is the
 485 size of weight vector \mathbf{w} . It is clear that the samples of the training set S
 486 should be greater than the number of neural network parameters, that is the
 487 dimension of the weight vector \mathbf{w} . Since the size of the neural network is
 488 small few training samples are required.

489 6.2. Recursive Implementation

490 The main drawback of the aforementioned approach is that a large num-
 491 ber of samples, as provided by the user's interaction through a set of relevant/
 492 irrelevant data is required for training the non-linear classifier. However, in
 493 real-life applications this large training set usually are constructed from data
 494 taken from quite different environmental conditions. This deteriorates the
 495 performance of the neural network model since, on the one hand, it contains

496 contradictory samples, and on the other, it is an averaging solution over
 497 data taken from different environmental conditions. To address the afore-
 498 mentioned difficulties a recursive implementation is proposed in this paper.
 499 This implementation requires few training samples that express the current
 500 user’s feedback through the selection of a set of relevant / irrelevant data.
 501 Then, the modification of the neural network weights is accomplished by
 502 minimizing the current network error with, however, a minimal modifica-
 503 tion of the previous neural network knowledge or in other words a minimum
 504 modification of the neural network weights.

505 In particular, let us assume that \mathbf{w}_b is the weights of the neural network
 506 classifier. Then, we denote as \mathbf{w}_n the new neural networks weights after
 507 the implementation of the recursive algorithm. Assuming that the new neu-
 508 ral network weights \mathbf{w}_n are related with the weights \mathbf{w}_b with a very small
 509 modification \mathbf{dw} , we can provide a very efficient training algorithm able to
 510 re-adjust the performance of the neural network classifier to the current user’s
 511 preferences.

512 In particular, let us denote that the output of the neural network classifier
 513 at a given time instance (e.g., the i^{th} frame) is $\mathbf{d}_i = f_{\mathbf{w}_b}(\mathbf{p}_i)$. In this case, we
 514 assume that the weights \mathbf{w}_b are used for the classifier. Let us now assume
 515 that the output of the classifier is not correct. The user provides the output
 516 target through user’s interaction. Let us denote this target as $\hat{\mathbf{d}}_i$. Variable
 517 $\hat{\mathbf{d}}_i$ refers to the actual output of the neural network classifier as provided by
 518 the user’s interaction. Then, the small amount of \mathbf{dw} is estimated through
 519 the following equation (see [49], [50])

$$\mathbf{e} = \hat{\mathbf{d}}_i - \mathbf{d}_i = \mathbf{a}^T \cdot \mathbf{dw} \quad (14)$$

520 Equation (14) is derived by using Taylor series expansion on the neu-
 521 ral network model (see [50]). In equation (14) vector \mathbf{a} contains elements
 522 of the previous network weights \mathbf{w}_b . Using only equation (14), we cannot
 523 estimate the small perturbation \mathbf{dw} . This is due to the fact that only one
 524 (or few) equations are not enough to estimate the multi-dimensional vector
 525 \mathbf{dw} . For this reason, an additional constraint is required. In this approach,
 526 the variable \mathbf{dw} is estimated through the minimum modification of the pre-
 527 vious network knowledge. Another approach is given by the minimization
 528 of the norm of the perturbation \mathbf{dw} resulting in the following minimization
 529 approach

$$\hat{\mathbf{d}}\mathbf{w} = \mathit{argmin}\|\mathbf{d}\mathbf{w}\| \quad (15)$$

subject to

$$\mathbf{e} = \hat{\mathbf{d}}_i - \mathbf{d}_i = \mathbf{a}^T \cdot \mathbf{d}\mathbf{w} \quad (16)$$

530 The previous equations express a convex minimization problem. There-
 531 fore, it can be easily solved either analytically using Lagrange multipliers or
 532 arithmetically using the gradient projection method [51]. In this paper, we
 533 adopt the second option to reduce computational complexity. This is due to
 534 the fact that the arithmetic approach is actually an iterative method. Thus,
 535 we can restrict the number of iterations with respect to the computational
 536 cost needed. In particular, the method starts from a feasible solution, i.e.,
 537 an arbitrary solution that satisfies the constraints. Then, this solution is
 538 iteratively updated according to the gradient of the square function $\|\mathbf{d}\mathbf{w}\|$,
 539 as being projected on the hyperplane defined by the constraint. It seems
 540 that this modification is very efficient for correcting erroneous mistakes of
 541 the HMM model combined with the particle filters (see the section of exper-
 542 imental results).

543 Summarizing, relevance feedback is a methodology of dynamically updat-
 544 ing the system response, by either modifying the system parameters or the
 545 entire system itself, through information provided by the user, regarding the
 546 relevance of a set of few samples selected by the user and feedback to the
 547 system. We have used a non-linear relationship, modeled through a dynamic
 548 neural network architecture, for the relevance feedback implementation. In
 549 particular, at specific or randomly selected time instances, the user interacts
 550 with the system, by indicating the perfect target output (real task) for this
 551 particular time instance (captured frame). This selected image frame is feed-
 552 back to the proposed adaptable architecture in order to improve the neural
 553 network response at possible future samples. In a nutshell, the proposed
 554 relevance feedback framework is described in the following:

- 555 • At a random image frame the user interacts with the system by pro-
 556 viding the perfect target output for this frame.
- 557 • The respective output of the HMMs (observation probabilities per HMM
 558 model for the selected image frame) are feedback to the adaptable neu-
 559 ral network architecture.

- 560 • The gradient project method is used to estimate the small perturbation
561 $d\mathbf{w}$ using equations (15), (16) and information of the sample given by
562 the user.
- 563 • The weights of the neural network are updated using the relation $\mathbf{w}_n =$
564 $\mathbf{w}_n + d\mathbf{w}$.
- 565 • The updated network is used to recalculate the task observation prob-
566 abilities.

567 7. Experiments and Results

568 We experimentally verified the applicability of the described methods. To
569 this end, we have acquired some very challenging videos from the production
570 line of a major automobile manufacturer (see [52]). Our previous efforts to
571 apply a detection-tracking scheme have failed due to the low resolution, the
572 heavy occlusions and the illumination changes (e.g., due to welding sparks).

573 7.1. Experimental Setup

574 The production cycle on the production line included tasks of picking
575 several parts from racks and placing them on a designated cell some meters
576 away, where welding took place. Each of the above tasks was regarded as
577 a class of behavioral patterns that had to be recognized. The information
578 acquired from this procedure can be used for the extraction of production
579 statistics or anomaly detection. Partial or total occlusions due to the racks
580 made the classification task difficult to effect.

581 The behaviors (tasks) we were aiming to model in the examined applica-
582 tion are briefly the following:

- 583 1. One worker picks part #1 from rack #1 and places it on the welding
584 cell.
- 585 2. Two workers pick part #2a from rack #2 and place it on the welding
586 cell.
- 587 3. Two workers pick part #2b from rack #3 and place it on the welding
588 cell.
- 589 4. A worker picks up parts #3a and #3b from rack #4 and places them
590 on the welding cell.
- 591 5. A worker picks up part #4 from rack #1 and places it on the welding
592 cell.

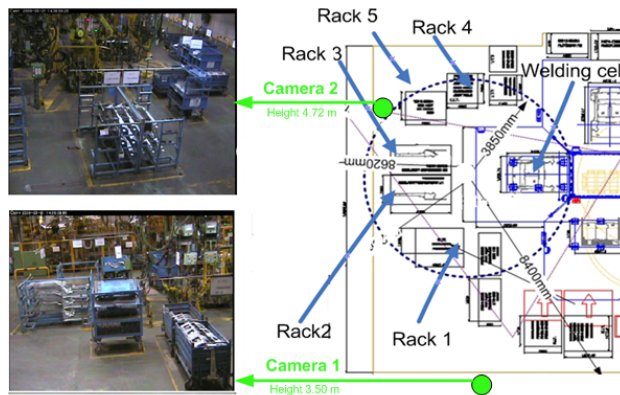


Figure 4: Depiction of a workcell along with the position our camera (camera 1) and the racks #1-5. The recognized behaviors are associated with transferring each part from the respective pallet and putting it on the welding cell.

593 6. Two workers pick up part #5 from rack #5 and place it on the welding
 594 cell.

595 In addition to these we had a null task (referenced as task 7), during
 596 which the workers were idle or absent.

597 The workspace configuration and the cameras' positioning is given in Fig.
 598 4. A sample task (task 2) is presented in Fig. 5. The work cycle that we
 599 sought to model, despite the noise and the several outliers (e.g., persons
 600 walking into the working cell, vehicles passing by etc), remains a structured
 601 process and is a good candidate to model with holistic features.

602 For our experiments, we have used 20 sequences representing full assem-
 603 bly cycles, each one containing each of the defined behaviors The length
 604 of each sequence ranges from 2000 frames to 4000¹. The annotation has
 605 been done manually. The dataset that we used is unique in the sense that it
 606 presents some well defined tasks which are executed in a repetitive and rather
 607 structured manner, providing several samples, which is good for learning (of
 608 course there is intra-class variability between the same tasks but still the
 609 resulting time series are correlated and can be learned and recognized). Fur-

¹We are going to make the dataset publicly available. It is currently available for review purposes on <http://www.4shared.com/dir/sYeCqK5d/SignalProcessingVideoAnalytics.html> (folder:dataset1 - password:xyz543)



Figure 5: Typical execution of task 2. The relatively low resolution and the several occlusions and self occlusions make very difficult the task of tracking thus necessitating a holistic method

610 furthermore, it includes phenomena such as occlusions, illumination changes
 611 and uniform appearance of humans, which make reliable tracking rather un-
 612 realistic, though perfectly suitable for approaches based on holistic features.
 613 Finally, partially overlapping views are available, which facilitates occlusions
 614 handling in our later research steps.

615 7.2. Holistic representation and online classification

616 To represent each video frame with a feature vector, we followed the
 617 method described in the subsection 4.1. For capturing the spatio-temporal
 618 variations we have set the parameters at $\varsigma = 10$ and $\tau = 70$. We have chosen
 619 to use the Zernike moments up to sixth order along with the center of gravity
 620 and the area, as feature vector. The higher the order of moments that we
 621 employ, the more detailed the region reconstruction will be, but also the more
 622 processing power will be required. Limiting the order of moments used is also
 623 justified by the fact that the details captured by higher order moments have
 624 much higher variability and are more sensitive to noise.

625 Specifically, we employed the complex moments $A_{00}, A_{11}, A_{20}, A_{22}, A_{31},$
 626 $A_{33}, A_{40}, A_{42}, A_{44}, A_{51}, A_{53}, A_{55}, A_{60}, A_{62}, A_{64}, A_{66}$, for each of which we
 627 used the norm and the angle, except for $A_{00}, A_{20}, A_{40}, A_{60}$, for which the angle

628 was always constant. Additionally the center of gravity and the area were
629 used, making a total of 31 parameters, thus providing an acceptable scene
630 reconstruction without a computationally prohibitive dimension. Zernike
631 moments have been calculated in rectangular regions of interest of approx-
632 imately 15000 pixels in each image, to limit the processing and allow real
633 time feature extraction.

634 For activity recognition we used three-state HMMs with one mixture
635 component per state to model each of the tasks described above, making a
636 discrete search space of size 7; this was a good trade-off between performance
637 and efficiency. In all cases, we employed full covariance matrices for the
638 adopted observation (mixture) models. We trained all our models using the
639 EM algorithm and we used the first ten scenarios for training and the rest
640 ten for testing.

641 Although the selection of the features to represent each frame is independ-
642 ent of the proposed classification method, we have compared the features
643 described in subsection 4.1 to the Local Motion Grid (LMG) features that
644 have been used on the same dataset (see our work [20], or [29]) to ensure high
645 accuracy. Using the same HMM configuration and a leave-one-out cross val-
646 idation scheme for 20 scenarios considering up to seven pre-segmented tasks
647 per scenario, we measured for cameras 1 and 2 a total accuracy of 53.57% and
648 67.14% for the LMG features [20], versus 84.14% and 86.42% respectively for
649 the PCH-Zernike representation [19]. This comparison justifies the sole use
650 of the PCH-Zernike features in the rest of the experiments.

651 We have compared the proposed method to some baseline methods the
652 first one being the standard HMM. We have taken sliding time windows of
653 constant size, which was equal to the mean duration of tasks in the training
654 set. For each of those windows we applied the HMM models that represent
655 each task and the winner was the one giving the higher likelihood. Using a
656 voting scheme we were able to classify each frame.

657 The second baseline method that we used for comparison was the echo
658 state network (ESN). We used a network of 1000 nodes, which was efficient
659 for real time execution and avoided overfitting. Increasing the number of
660 nodes would result in very high memory requirements without actual benefit
661 in accuracy or recall. It also had seven output nodes, each one of them
662 corresponding to a predicted task. The median of the last 101 estimations
663 was taken to ensure lower jitter in the output. We have used the Matlab
664 toolbox provided by the authors [27] using spectral radius 0.60, input scaling
665 0.3 and noise level 0.0003 after some experimentation for optimal results.

666 For our method we used only 100 particles, which was a good trade-
667 off, and we were able to perform the whole processing at a rate of about
668 20Hz, the most costly of which was the feature extraction. The confusion
669 matrix per task for a typical case is given in Fig. 6a. The learning phase
670 included learning the task durations using a Gaussian mixture model, the
671 task trajectories using a decision tree and the task models using HMMs.

672 For all methods we extracted recall and precision. Recall indicates the
673 number of true positives divided by the total number of positives in the
674 ground truth ($REC = TP/(TP + FN)$). Precision is the number of true posi-
675 tives divided by the number of true and false positives ($PRC = TP/(TP +$
676 $FP)$). The average precision and recall per task and the standard deviation
677 are given in Fig. 7, after performing twenty iterations. The overall results
678 are given in Table 1 and the confusion matrices are given in Fig. 6.

679 As expected the conventional HMM has the worst performance, mainly
680 due to the fact that the task durations may vary, while it uses a sliding
681 window of constant size. Although the ESN performs generally well, each
682 task may be mistaken for almost any other one, as becomes clear from the
683 confusion matrix. This happens mainly because actually only the most recent
684 observations guide the prediction, thus ignoring the whole history (effective
685 Markovian behavior [28]). Providing constraints given the history of tasks
686 helps to discard erroneous task transitions by utilizing a decision tree, as
687 we have explained in subsection 5.1. This becomes obvious when observing
688 the confusion matrix corresponding to our method, where errors seem to be
689 distributed among mutually accessible tasks. Clearly, the more restrictive
690 the structure of the tasks, the more effective such a scheme will be, because
691 the particles will be scattered around the most probable tasks according to
692 the observation history.

693 In our method the particle that was able to explain best the sequence
694 according to (3) was considered to be the winner. In all cases the work cycle,
695 which consisted of all tasks 1 to 6 and the null was successfully recognized.
696 In few cases the tasks were identified but were misaligned with the real ones;
697 this was mainly due to features' imperfections, which sometimes gave rise to
698 the wrong tasks, due to occlusions and noise.

699 7.3. User feedback

700 Regarding the user feedback mechanism, we firstly used a feed-forward
701 neural network model for estimating the distribution probability of a frame to

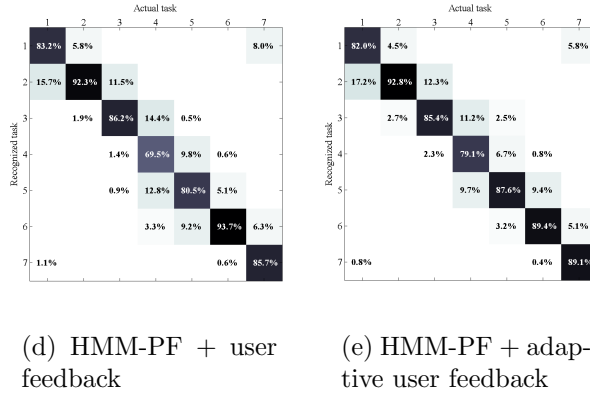
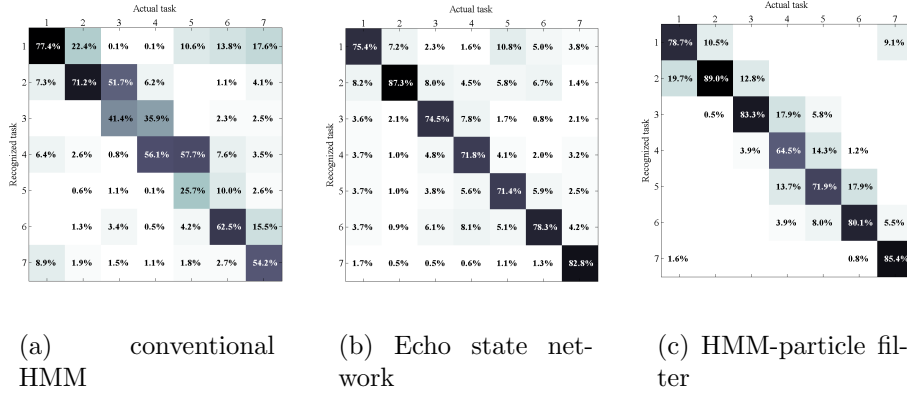


Figure 6: Confusion matrices for the various methods

702 belong to one of the seven available categories. The feed-forward neural net-
 703 work model (see Fig. 8) has one hidden layer. In our experiments 15 hidden
 704 neurons have been selected. As the number of hidden neurons increases the
 705 complexity of the neural network training significantly increases as well and
 706 generalization performance of the network decreases. In our simulations, we
 707 have seven output neurons for the neural model that indicate the probability
 708 for one of the seven available categories. As input layer of the neural net-
 709 work model we have the seven observation probabilities of the current frame
 710 produced by HMM. The neural network model re-adjusted the probabilities
 711 according to the knowledge provided through supervised training. We used
 712 the scenarios 1 to 3 to map the incoming maximum observation probabilities
 713 given each task (as provided by the HMM) to the ideal ones.

Table 1: Overall precision and recall for 10 test scenarios.

Method	HMM	ESN	HMM-PF	HMM-NN	HMM-NN recurs.
precision	0.603	0.777	0.797	0.851	0.871
recall	0.565	0.772	0.788	0.846	0.863

Table 2: Overall precision and recall for HMM-PF rectified by the recursive neural network.

number of particles	30	60	100	200
precision (%)	87.52	87.88	87.12	87.45
recall (%)	86.63	86.83	86.28	86.44

714 The neural network output re-adjusts the probabilities of the combined
715 HMM + particle filter model in order to increase the efficiency of the work-
716 flow recognition module. There are, however, several cases, where the per-
717 formance of the neural network model is not satisfactory. For this reason,
718 we used the on-line learning mechanism for dynamically updating the pa-
719 rameters of the neural network model in order to satisfy with more efficiency
720 the target outputs. In particular, at random selected time instances the user
721 provided the perfect target outputs by setting the probability of the desired
722 task to 1 and the other ones close to zero. We did that for 70 random sam-
723 ples. Then, the system updated the parameters of the neural network model
724 so that i) the corrected target output was satisfied as much as possible, while
725 simultaneously ii) the minimum modification of the previous network weights
726 (that is previous neural network performance) was satisfied. Then, the new
727 estimated weights were used for predicting the workflow state of future image
728 frames.

729 Clearly the performance was improved compared to the previous approach
730 in terms of precision and recall using the same number of particles (see Fig.
731 6e, Fig. 7 and Table 1). The high performance verifies that the network is
732 able to adapt its weights to minimize the error according to the most recent
733 input samples.

734 When using the recursive neural network we discovered that the increase
735 of the particles does not significantly affect the overall precision and recall.
736 The results are presented in Table 2 given for 30, 60, 100 and 200 particles.
737 This implies that the proposed feedback scheme provided a good estimate
738 even with low number of samples.

739 **8. Conclusions**

740 In this work we have proposed a novel online framework for behavior
741 recognition in workflows in real-time. In the context of the framework we
742 have handled two important problems for behavior recognition: (a) online be-
743 havior classification through a Bayesian filter which is approximated through
744 particles driven by HMM and (b) rectification of erroneous classifications
745 through interaction with the user.

746 The holistic features gave a good scene representation, thus helping us
747 bypass the difficult tasks of detection and tracking that fail in such com-
748 plex sequences. The conventional application of Viterbi to obtain optimal
749 result would make the recognition task infeasible given the fact that no ini-
750 tial and end sequence points were known. Furthermore, our method did not
751 rely on the Markovian assumption, which is not appropriate for monitoring
752 workflows.

753 The proposed methods have been applied with promising results in some
754 very challenging real sequences from an automobile manufacturing process.
755 The good online recognition rates achieved by the particle filter/HMM method
756 are additionally improved significantly when we employed the neural network
757 based rectification scheme that incorporated user feedback. The recursive
758 scheme seemed to perform even better and required fewer particles to achieve
759 similar performance.

760 In the case of very long tasks, it would be meaningful to have more parti-
761 cles and maintain more hypotheses, e.g., in the case of consecutive workflows.
762 If the workflows can be separated, e.g., when for several particles all expected
763 tasks are considered finished, it would be more practical to reset the particles
764 and to start from the beginning.

765 In our experiments we used unique tasks, which always had to appear due
766 to the industrial assembly workflow requirements. In different settings, where
767 probably omission or repetition of tasks would be possible we would only need
768 to model these omissions/repetitions as prior knowledge, i.e., possible paths
769 in the tree that we defined in subsection 5.1.2. What is only needed is a
770 good estimation of which task is probable to appear next; that requirement
771 is covered by the proposed representation.

772 In our future work we are going to investigate some less structured sce-
773 narios with more complex interactions and tasks with even higher variability,
774 also considering different viewpoints via fusion schemes.

775 **References**

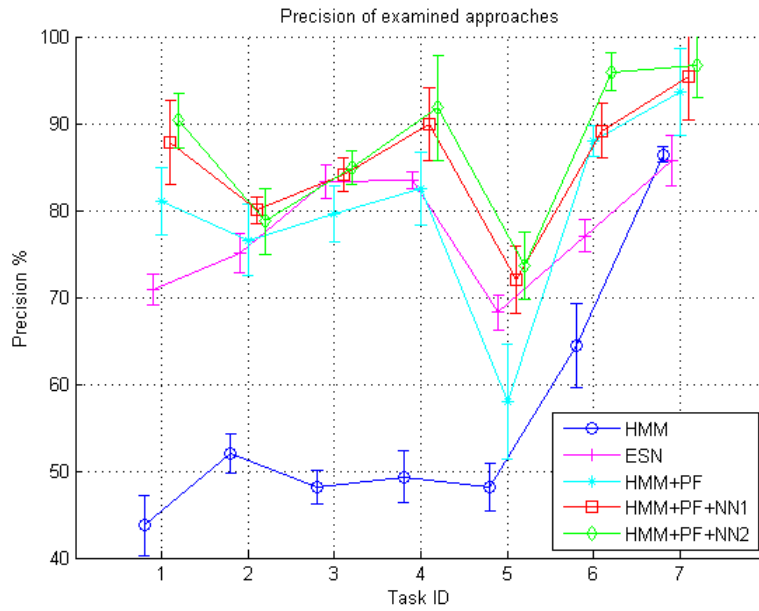
- 776 [1] D. Kosmopoulos, A. Voulodimos, T. Varvarigou, Robust human behav-
777 ior modeling from multiple cameras, in: Pattern Recognition (ICPR),
778 2010 20th 697 International Conference on, 2010, pp. 3575–3578.
- 779 [2] N. Doulamis, A. Voulodimos, D. Kosmopoulos, T. Varvarigou, Enhanced
780 human behavior recognition using hmm and evaluative rectification, in:
781 ACM Multimedia, ARTEMIS Workshop, 2010.
- 782 [3] Y. Wang, G. Mori, Human action recognition by semilattent topic mod-
783 els, Pattern Analysis and Machine Intelligence, IEEE Transactions on
784 31 (10) (2009) 1762 –1774.
- 785 [4] Q. Hu, L. Qin, Q. Huang, S. Jiang, Q. Tian, Action recognition using
786 spatial-temporal context, in: Pattern Recognition (ICPR), 2010 20th
787 International Conference on, 2010, pp. 1521 –1524.
- 788 [5] S. Ali, M. Shah, Human action recognition in videos using kinematic
789 features and multiple instance learning, Pattern Analysis and Machine
790 Intelligence, IEEE Transactions on 32 (2) (2010) 288 –303.
- 791 [6] A. Yao, J. Gall, L. Van Gool, A hough transform-based voting frame-
792 work for action recognition, in: Computer Vision and Pattern Recogni-
793 tion (CVPR), 2010 IEEE Conference on, 2010, pp. 2061 –2068.
- 794 [7] T. Wada, T. Matsuyama, Multiobject behavior recognition by event
795 driven selective attention method, Pattern Analysis and Machine Intel-
796 ligence, IEEE Transactions on 22 (8) (2000) 873 –887.
- 797 [8] D. Sadlier, N. O’Connor, Event detection in field sports video using
798 audio-visual features and a support vector machine, Circuits and Sys-
799 tems for Video Technology, IEEE Transactions on 15 (10) (2005) 1225
800 – 1233.
- 801 [9] M.-H. Hung, C.-H. Hsieh, Event detection of broadcast baseball videos,
802 Circuits and Systems for Video Technology, IEEE Transactions on
803 18 (12) (2008) 1713 –1726.
- 804 [10] I. Laptev, P. Perez, Retrieving actions in movies, in: Computer Vision,
805 2007. ICCV 2007. IEEE 11th International Conference on, 2007, pp. 1
806 –8.

- 807 [11] A. Bobick, A. Wilson, A state-based approach to the representation and
808 recognition of gesture, *Pattern Analysis and Machine Intelligence*, IEEE
809 *Transactions on* 19 (12) (1997) 1325–1337.
- 810 [12] M.-H. Yang, N. Ahuja, Extraction and classification of visual motion
811 patterns for hand gesture recognition, in: *Computer Vision and Pattern*
812 *Recognition*, 1998. Proceedings. 1998 IEEE Computer Society Confer-
813 ence on, 1998, pp. 892–897.
- 814 [13] R. Hamid, S. Maddi, A. Bobick, M. Essa, Structure from statistics -
815 unsupervised activity analysis using suffix trees, in: *Computer Vision*,
816 2007. ICCV 2007. IEEE 11th International Conference on, 2007, pp. 1
817 –8.
- 818 [14] R. Poppe, A survey on vision-based human action recognition, *Image*
819 *and Vision Computing* 28 (6) (2010) 976–990.
- 820 [15] W. Hu, T. Tan, L. Wang, S. Maybank, A survey on visual surveillance
821 of object motion and behaviors, *Systems, Man and Cybernetics, Part*
822 *C*, *IEEE Transactions on* 34 (3) (2004) 334–352.
- 823 [16] L. R. Rabiner, A tutorial on hidden Markov models and selected ap-
824 plications in speech recognition, *Proceedings of the IEEE* 77 (2) (1989)
825 257–286.
- 826 [17] S. P. Chatzis, D. I. Kosmopoulos, T. A. Varvarigou, Robust sequential
827 data modeling using an outlier tolerant hidden Markov model, *IEEE*
828 *Transactions on Pattern Analysis and Machine Intelligence* 31 (9) (2009)
829 1657–1669.
- 830 [18] Z. Zeng, J. Tu, B. Pianfetti, T. Huang, Audio-visual affective expression
831 recognition through multistream fused HMM, *IEEE Trans. Multimedia*
832 10 (4) (2008) 570–577.
- 833 [19] D. Kosmopoulos, S. Chatzis, Robust visual behavior recognition, *Signal*
834 *Processing Magazine*, *IEEE* 27 (5) (2010) 34–45.
- 835 [20] A. Voulodimos, H. Grabner, D. I. Kosmopoulos, L. J. V. Gool, T. A. Var-
836 varigou, Robust workflow recognition using holistic features and outlier-
837 tolerant fused hidden markov models, in: *ICANN* (1), 2010, pp. 551–560.

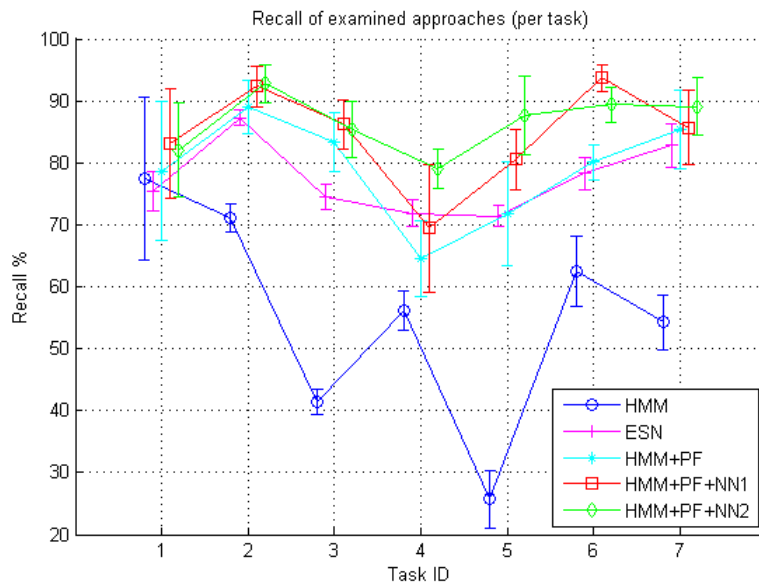
- 838 [21] S. Eickeler, A. Kosmala, G. Rigoll, Hidden markov model based con-
839 tinuous online gesture recognition, in: In Int. Conference on Pattern
840 Recognition (ICPR, 1998, pp. 1206–1208.
- 841 [22] F. Lv, R. Nevatia, Recognition and segmentation of 3-d human action
842 using hmm and multi-class adaboost, in: ECCV06, 2006, pp. IV: 359–
843 372.
- 844 [23] S. Fine, Y. Singer, N. Tishby, The hierarchical hidden markov model:
845 Analysis and applications, *Machine Learning* 32 (1) (1998) 41–62.
- 846 [24] N. Padoy, D. Mateus, D. Weinland, M.-O. Berger, N. Navab, Work-
847 flow Monitoring based on 3D Motion Features, in: Workshop on Video-
848 Oriented Object and Event Classification in Conjunction with ICCV
849 2009, IEEE, Kyoto Japan, 2009, pp. 585–592.
- 850 [25] N. Oliver, A. Garg, E. Horvitz, Layered representations for learning and
851 inferring office activity from multiple sensory channels, *Comput. Vis.*
852 *Image Underst.* 96 (2) (2004) 163–180.
- 853 [26] T. Xiang, S. Gong, Optimising dynamic graphical models for video con-
854 tent analysis, *Comput. Vis. Image Underst.* 112 (2008) 310–323.
- 855 [27] H. Jaeger, W. Maass, J. Principe, Special issue on echo state networks
856 and liquid state machines, *Neural Networks* 20 (3) (2007) 287 – 289.
- 857 [28] C. Gallicchio, A. Micheli, Architectural and markovian factors of echo
858 state networks, *Neural Networks* 24 (5) (2011) 440 – 456.
- 859 [29] G. V. Veres, H. Grabner, L. Middleton, L. J. V. Gool, Automatic work-
860 flow monitoring in industrial environments, in: ACCV (1), 2010, pp.
861 200–213.
- 862 [30] X. Zhang, G.-Y. Xu, X.-L. Xiao, L.-M. Tao, Online analysis of hierar-
863 chical events in meetings, in: Usability and Internationalization, HCI
864 and Culture, Vol. 2, 2007, pp. 472 –479.
- 865 [31] X. Xiaoling, L. Layuan, Real time analysis of situation events for in-
866 telligent surveillance, in: Computational Intelligence and Design, 2008.
867 ISCID '08. International Symposium on, Vol. 2, 2008, pp. 122 –125.

- 868 [32] D. Zhang, X. Ning, X. Liu, Smc method for online prediction in hidden
869 markov models, *Kybernetes* 38 (10) (2009) 1819–1827.
- 870 [33] P. Fearnhead, P. Clifford, On-line inference for hidden markov models
871 via particle filters, *Journal Of The Royal Statistical Society Series B*
872 65 (4) (2003) 887–899.
- 873 [34] N. Doulamis, A. Doulamis, Evaluation of relevance feedback schemes in
874 content-based retrieval systems, *Signal Processing: Image Communica-*
875 *tion* 21 (4) (2006) 334 – 357.
- 876 [35] J. J. Rocchio, Relevance feedback in information retrieval, in: G. Salton
877 (Ed.), *The Smart retrieval system - experiments in automatic document*
878 *processing*, Englewood Cliffs, NJ: Prentice-Hall, 1971, pp. 313–323.
- 879 [36] N. D. Doulamis, A. D. Doulamis, Evaluation of relevance feedback
880 schemes in content-based in retrieval systems, in: *Signal process-*
881 *ing:image communication*, Vol. 21, Elsevier, 2006, pp. 334–357.
- 882 [37] A. Oerlemans, J. T. Rijsdam, M. S. Lew, Real-time object tracking
883 with relevance feedback, in: *Proceedings of the 6th ACM international*
884 *conference on Image and video retrieval, CIVR '07*, ACM, New York,
885 NY, USA, 2007, pp. 101–104.
- 886 [38] Z. Chengcui, C. Wei-Bang, C. Xin, Y. Lin, J. John, A multiple instance
887 learning and relevance feedback framework for retrieving abnormal in-
888 cidents in surveillance videos, *Journal of Multimedia* 5 (2010) 310–321.
- 889 [39] A. Doulamis, N. Doulamis, Generalized nonlinear relevance feedback
890 for interactive content-based retrieval and organization, *IEEE Trans. on*
891 *Circuits and Systems for Video Technology* 14 (5) (2004) 656–671.
- 892 [40] X. S. Zhou, T. S. Huang, Relevance feedback in image retrieval: A
893 comprehensive review, *Multimedia Systems* 8 (6) (2003) 536–544.
- 894 [41] A. Lam, A. K. Roy-Chowdhury, C. R. Shelton, Interactive event search
895 through transfer learning (2011) 157–170.
- 896 [42] X. Sun, M. Chen, A. Hauptmann, Action recognition via local descrip-
897 tors and holistic features, in: *IEEE Conference on Computer Vision and*
898 *Pattern Recognition*, 2009, pp. 58–65.

- 899 [43] T. Xiang, S. Gong, Beyond tracking: modelling activity and under-
900 standing behaviour, *International Journal of Computer Vision* 67 (2006)
901 21–51.
- 902 [44] M. S. Arulampalam, S. Maskell, N. Gordon, A tutorial on particle filters
903 for online nonlinear/non-gaussian bayesian tracking, *IEEE Transactions*
904 *on Signal Processing* 50 (2) (2002) 174–188.
- 905 [45] D. Arnaud, G. Simon, A. Christophe, On sequential monte carlo sam-
906 pling methods for bayesian filtering, *Statistics and Computing* 10 (3)
907 (2000) 197–208.
- 908 [46] G. Seber, C. Wild, *Nonlinear Regression*, Wiley, Hoboken, New Jersey,
909 2003.
- 910 [47] L. Breiman, Random forests, *Mach. Learn.* 45 (2001) 5–32.
- 911 [48] D. E. Rumelhart, G. E. Hinton, R. J. Williams, *Learning representations*
912 *by back-propagating errors*, MIT Press, Cambridge, MA, USA, 1988, pp.
913 696–699.
- 914 [49] A. Doulamis, Adaptable neural networks for objects’ tracking re-
915 initialization, Vol. 5769 LNCS of *Lecture Notes in Computer Science*,
916 2009, pp. 715–724.
- 917 [50] A. Doulamis, Knowledge extraction in stereo video sequences using
918 adaptive neural networks, *Intelligent Multimedia Processing with Soft*
919 *Computing*, Springer-Verlag Berlin Heidelberg, 2005, pp. 235–252.
- 920 [51] D. G. Luenberger, *Linear and nonlinear programming*, Addison-Wesley,
921 1984.
- 922 [52] A. Voulodimos, D. Kosmopoulos, G. Vasileiou, E. Sardis, A. Doulamis,
923 V. Anagnostopoulos, C. Lalos, T. Varvarigou, A dataset for workflow
924 recognition in industrial scenes, in: *IEEE Int. Conference on Image*
925 *Processing*, 2011.



(a) Precision



(b) Recall

34
 Figure 7: Comparison of mean precision-recall metrics and standard deviation for conventional HMM, ESN, HMM-PF, HMM-PF-NN1 (non-linear), HMM-PF-NN2 (recursive)

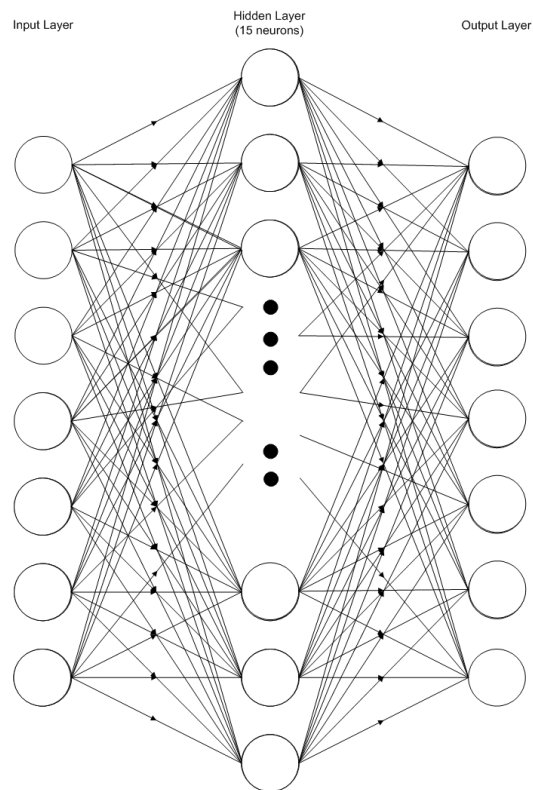


Figure 8: The feed-forward neural network used in the proposed rectification scheme. The input and output layers consist of seven nodes each (one for each task) and the hidden layer comprises 15 nodes