

1 **Citation:**

2 *E. Bicho, W. Erlhagen, L. Louro, E. Costa e Silva, R. Silva, N. Hipolito, "A*
3 *dynamic field approach to goal inference, error detection and anticipatory action*
4 *selection in human-robot collaboration", In Dautenhah, Kerstin and Saunders,*
5 *Joe (Eds), New Frontiers in Human Robots Interaction, [Advances in Interaction*
6 *Studies, 2] 2011. Vi, 332 (pp. 135-164, John Benjamins Publishing Company.*
7 *Hardbound: ISBN 978 90 272 0455 4 eBook: ISBN 978 90 272 8339 9*

8

9 **Running head:** A DNF approach to human-robot joint action

10

11 **A dynamic field approach to goal inference, error**
12 **detection and anticipatory action selection in human-**
13 **robot collaboration**

14 Estela Bicho^{a1}, Wolfram Erlhagen^b, Luís Louro^a, Eliana Costa e Silva^a, Rui
15 Silva^a, Nzoji Hipolito^a

16

^aDepartment of Industrial Electronics

17

^bDepartment of Mathematics and Applications

18

University of Minho, Portugal

19

20 **Abstract**

21 In this chapter we present results of our ongoing research on efficient and
22 fluent human-robot collaboration that is heavily inspired by recent
23 experimental findings about the neurocognitive mechanisms supporting
24 joint action in humans. The robot control architecture implements the joint
25 coordination of actions and goals as a dynamic process that integrates
26 contextual cues, shared task knowledge and the predicted outcome of the
27 user's motor behavior. The architecture is formalized as a coupled system of
28 dynamic neural fields representing a distributed network of local but
29 connected neural populations with specific functionalities. We validate the
30 approach in a task in which a robot and a human user jointly construct a toy
31 'vehicle'. We show that the context-dependent mapping from action

¹Email addresses: estela.bicho@dei.uminho.pt (corresponding author),
wolfram.erlhagen@math.uminho.pt (Wolfram Erlhagen), louro@dei.uminho.pt (Luís
Louro), esilva@dei.uminho.pt (Eliana Costa e Silva), rsilva@dei.uminho.pt (Rui Silva) and
nhipolito@dei.uminho.pt (Nzoji Hipolito).

32 observation onto appropriate complementary actions allows the robot to
33 cope with dynamically changing joint action situations. More specifically,
34 the results illustrate crucial cognitive capacities for efficient and successful
35 human-robot collaboration such as goal inference, error detection and
36 anticipatory action selection.

37

38 **1. Introduction**

39

40 As robot systems are moving as assistants into human everyday life, the
41 question how to design robots capable of acting as sociable partners in
42 collaborative joint activity becomes increasingly important (Breazeal, 2004;
43 Fong, Nourbakhsh and Dautenhahn, 2003). Useful and efficient human-
44 robot collaboration requires that both teammates coordinate and synchronize
45 their actions and decisions in a shared task. In order to decrease the
46 workload of the human and to increase user satisfaction, the robot should
47 equally contribute to this coordination effort. This necessarily means that
48 the robot should be endowed with cognitive capacities such as action
49 understanding, action monitoring and goal inference. Humans achieve their
50 remarkable fluent organization of joint action by anticipating the motor
51 intentions of others (Sebanz, Bekkering and Knoblich, 2006). In our
52 everyday social interactions we continuously monitor the actions of our
53 partners, interpret them effortlessly in terms of their outcomes and use these
54 predictions to select adequate complementary behaviours. Very often this
55 happens without the need for explicit verbal communication. Imagine for
56 the instance the joint action task of preparing a dinner table. The way that a
57 partner grasps a certain object, e.g., a coffee cup, transmits to the observer
58 important information about the ultimate goal of the action. Depending on
59 the grip type, the partner may want to place the cup on the table or,
60 alternatively, has the intention to hand it over to the co-actor. Being able to
61 predict the goal of the whole action sequence at the time of the grasping
62 allows the observer to timely prepare for receiving the cup, or to initiate the
63 selection of another object for the dinner table. However, even in routine
64 joint activity the co-actor may perform actions that are in some way
65 incorrect or inappropriate. The partner may for instance want to hand over a
66 second spoon for the sugar bowl or may pick a cup without having already
67 placed the saucer on the table. Being able to evaluate the predicted
68 outcomes of the co-actor's actions with respect to the (sub)goals of the
69 shared task is thus a fundamental capacity for efficient and successful joint
70 action. It allows the observer to overrule a familiar response (e.g., accepting
71 the spoon) or to initiate an adequate corrective behaviour (e.g., quickly
72 grasping a saucer for placing it on the table).

73 This chapter presents results of our ongoing research towards creating
74 socially intelligent robots that are able to flexibly adjust their goal-directed
75 behaviours depending on the predicted outcomes of actions of their human
76 partners. For the experiments we used a more complex version of a joint
77 assembly task introduced in our previous work (Bicho et al., 2009; Bicho et
78 al., 2010) in which the human-robot team has to assemble a toy object from
79 its components. The focus is on implementing and testing in the robot
80 context-sensitive action monitoring and anticipatory action selection
81 capacities. Our approach is heavily inspired by recent experimental and
82 theoretical findings about the neurocognitive mechanisms underlying joint
83 action in humans (Bekkering et al., 2009; Sebanz, Bekkering and Knoblich,
84 2006). We believe that designing cognitive control architectures on the basis
85 of these mechanisms defines a very promising research direction to reduce
86 the significant imbalance in social and cognitive skills between human and
87 robot that still exists today. Ultimately, implementing a human-like joint
88 action model in the robot will contribute to more natural human-robot
89 collaboration since the teammates will become more predictable to each
90 other. This in turn will increase the acceptance by humans.

91 An impressive body of experimental evidence from behavioural and
92 neurophysiological studies investigating action and perception in a social
93 context shows that when we observe others' actions, corresponding motor
94 representations in our motor system become activated (for reviews see
95 (Wilson and Knoblich, 2005; Rizzolatti and Craighero, 2004). These
96 findings have been interpreted as supporting the hypothesis that the mere
97 perception of actions automatically increases the likelihood of the
98 performance of those actions, without the person's conscious awareness.
99 During the last decade, the idea of an obligatory and direct perception-
100 action link has inspired robotics work mainly in the domain of learning by
101 imitation and social development (Billard et al., 2008; Erlhagen et al., 2006;
102 Demiris and Johnson, 2003; Alissandrakis et al., 2002; Schaal, 1999). A
103 major insight of this work is that the matching between action observation
104 and action execution has to solve the correspondence problems that exist
105 between agents with dissimilar embodiment. Moreover, the metrics of the
106 mapping should be highly task-dependent and may range from the level of
107 movement kinematics to the level of desired end states or action goals.

108 While an automatic facilitation of corresponding motor representations
109 during action observation may support social learning, it is normally not
110 beneficial for cooperative joint action tasks that require the facilitation of
111 complementary motor programs. An alternative proposal for a functional
112 role of the automatic action resonance mechanism suggests that it
113 contributes to understanding the actions of other individuals during social
114 interactions (Rizzolatti and Craighero, 2004; Wilson and Knoblich, 2005).

115 The key idea is that the observer performs an internal motor simulation to
116 predict the consequences of perceived actions using knowledge of his or her
117 actions and motor intentions. During joint action, the representation of the
118 inferred goal of the co-actor together with representations of prior task
119 knowledge may then automatically bias the observer's decision process
120 towards selecting an adequate complementary behaviour. In line with this
121 hypothesis, the findings of a recent behavioural study suggest that the
122 perception-action coupling appears to be indeed to some extent under the
123 control of task and goal representations. By comparing motor planning in an
124 imitative and a cooperative setting van Schie et al. (2008) demonstrated the
125 reversal of the automatic action congruency effect. In the cooperative
126 setting, people were faster to respond to an observed action with a
127 complementary behaviour compared to the matching behaviour.

128 The robot control architecture for human-robot collaboration implements
129 such a context-sensitive, i.e. flexible, mapping between action observation
130 and action execution. The coordination of actions and decisions among the
131 teammates is modeled as a dynamic process that builds on the continuous
132 integration of input from representations of the inferred goal of observed
133 actions (obtained through motor simulation), contextual cues (e.g., location
134 of objects in the scene) and shared task knowledge (e.g., assembly plan).
135 The representation of the complementary action that gets the strongest
136 support will win the dynamic competition process among all possible
137 complementary behaviours. As a theoretical framework we have used the
138 Dynamic Neural Field (DNF) approach to robotics (Erlhagen and Bicho,
139 2006). Originally introduced as a simplified mathematical model for pattern
140 formation in neural populations (Amari, 1977; Wilson and Cowan, 1973),
141 DNFs have been later generalized and applied to the cognitive domain (for a
142 recent review see Schöner, 2008). The architecture of DNFs reflects the
143 hypothesis that strong recurrent interactions in local populations of neurons
144 form a basic mechanism of cortical information processing. These
145 interactions support the existence of self-stabilized representations that
146 allow the cognitive agent for instance to compensate for temporally missing
147 sensory input, or to anticipate future environmental inputs that may inform
148 the decision about a specific goal-directed behaviour.

149 The DNF-model of joint action forms a complex dynamical system
150 consisting of a distributed network of reciprocally connected neural
151 populations that integrate and represent in their activation patterns task-
152 relevant information. For the experimental validation of the model in the
153 joint construction task we assume that both agents share the knowledge of
154 the assembly plan representing the sequential execution of subgoals. Since
155 the construction work cannot be performed alone, each agent has to
156 continuously monitor and evaluate the co-actor's actions in order to

157 guarantee success. The results reported here extend our previous work to a
158 more realistic and complex joint action context which includes situations in
159 which the co-actor's behaviour is only partially observable (due to occluding
160 surfaces) and the robot has to select among several possible complementary
161 actions. The focus of the study is on the dynamic interactions within the
162 DNF network that support the selection of an appropriate action in
163 anticipation of the co-actor's current goal. The anticipatory action control
164 includes situations in which the predicted effect of the observed action is
165 inconsistent with an efficient team performance and thus requires a
166 corrective response. The timely decision for such a response is possible
167 since the action planning process integrates continuously in time the activity
168 from connected populations representing a mismatch between the inferred
169 goal and the desired action effect in a specific joint action context.

170

171 The chapter is organized as follows: Section 2 introduces the joint
172 construction task and the robotic platform. Section 3 gives an overview
173 about the cognitive control architecture. Section 4 presents the basic
174 concepts of the dynamic field framework. The results of the human-robot
175 interactions are described in section 5. The chapter ends with a discussion of
176 concepts and results in section 6 and conclusions and outlook in section 7.

177

178

179 **2. Joint construction task**

180

181 -----Insert Figure 1 around here -----

182

183 To test the dynamic field architecture for human-robot collaboration we
184 have chosen the joint construction of a toy 'vehicle' from components that
185 are initially distributed on a table (see Figure 1). The toy object consists of
186 a round platform with an axle on which two wheels have to be attached and
187 each fixed with a nut. Subsequently, 4 different columns have to be plugged
188 into specific holes in the platform. The placing of another round object on
189 top of the columns finishes the task. The components were designed to limit
190 the workload for the vision and the motor system of the robot. Thus, the task
191 is completely symmetric in that both the human and the robot can make
192 assembly actions. It is assumed that each teammate is responsible to
193 assemble one side of the toy. Since the working areas of the human and the
194 robot do not overlap, the spatial distribution of components on the table
195 obliges the team to coordinate handing-over sequences. In addition, some
196 assembly steps require that one actor helps the other by holding still a part
197 in a certain position. It is further assumed that both partners know the

198 construction plan and keep track of the subtasks which have been already
199 completed by the team. The prior knowledge about the sequential execution
200 of the assembly work is represented in the DNF-architecture by pre-defined
201 connections between populations encoding subsequent assembly steps.
202 Since the desired end state does not uniquely define the logical order of the
203 construction, at each stage of the construction the execution of several
204 subtasks may be simultaneously possible.

205 The main challenge for the team is thus to efficiently coordinate in space
206 and time the decision about actions to be performed by each of the
207 teammates. The task is complex enough to show the impact of goal
208 inference, action understanding and action monitoring and evaluation on
209 complementary action selection.

210 The robot (ARoS) used in the experiments has been built in our lab (Silva,
211 Bicho and Erlhagen, 2008). It consists of a stationary torus on which a 7
212 DOFs AMTEC arm (Schunk GmbH) with a 3-fingers dexterous gripper
213 (Barrett Technology Inc.) and a stereo camera head are mounted. A speech
214 synthesizer (Microsoft Speech SDK 5.1) allows the robot to communicate
215 the result of its reasoning to the human user. For the control of the arm-hand
216 system we applied a global planning method in posture space that allows us
217 to integrate optimization principles derived from experiments with humans
218 (Costa et Silva et al, submitted). The goal is to guarantee robot motion that
219 is perceived by the human user as smooth and goal-directed.

220 The information about object class, position and pose is provided by the
221 vision system. The object recognition combines color-based segmentation
222 with template matching derived from earlier learning examples (Westphal et
223 al., 2008). The same technique is also used for the classification of object-
224 directed, static hand postures such as grasping and communicative gestures
225 such pointing and demanding an object.

226

227

228

229 **3. Cognitive architecture for joint action**

230

231 -----Insert Figure 2 around here -----

232

233 Figure 2 presents a sketch of the multi-layered robot control architecture. It
234 reflects neurocognitive mechanisms that are believed to support human joint
235 action (Bekkering et al., 2009). Each layer contains several neural
236 populations encoding information relevant for the joint assembly task.
237 Every population can receive input from multiple connected populations
238 that may be located in different layers. Rather than describing in detail the
239 schema of the hand-coded connections for the concrete assembly task (for

240 an example see the supplementary material in Bicho, Louro and Erlhagen,
241 2010) we give here an overview about the functional role of the different
242 layers and discuss the flow of information between layers with respect to
243 experimental findings that have inspired our work.

244 Ultimately, the architecture implements a context-dependent mapping
245 between observed action and executed action (Poljac, van Schie and
246 Bekkering, 2009; van Schie, Waterschoot and Bekkering, 2008; Erlhagen,
247 Mukovski and Bicho, 2006). The fundamental idea is that the mapping
248 takes place on the level of abstract motor primitives defined as whole
249 object-directed motor acts like reaching, grasping, placing, attaching or
250 plugging. These primitives encode the motor act in terms of an observable
251 end state or goal rather than in terms of a detailed description of the
252 movement kinematics (Rizzolatti and Craighero, 2004; Schaal, 1999). An
253 observed hand movement that is recognized by the vision system as a
254 particular primitive (e.g., top grip or side grip) is represented in the action
255 observation layer (AOL). The action simulation layer (ASL) implements the
256 idea that by automatically matching the co-actor's action onto its own
257 sensorimotor representations without executing it, the robot may simulate
258 the ongoing action and its consequences. The neural populations in ASL
259 represent entire chains of action primitives that are in the motor repertoire of
260 the robot (e.g., reaching-grasping-placing/plugging or reaching-grasping-
261 holding out). These chains are linked to representations of specific goals or
262 end states (e.g., attach wheel to base) which are represented by populations
263 in the intention layer (IL). Here the action chains are pre-coded, but in our
264 previous work we have addressed how they may autonomously develop
265 (Erlhagen, Mukovski and Bicho, 2006; Erlhagen et al., 2007). This chained
266 organization of motor primitives is motivated by recent findings of specific
267 neural populations in the inferior parietal lobe of monkey which is known to
268 be part of the matching system in monkey. Fogassi and colleagues (2005)
269 described neurons that fire at the time of a specific motor act (e.g., grasping)
270 in dependence of the ultimate goal of the action sequence in which the act is
271 embedded (e.g., grasping for placing versus grasping for eating). If a chain
272 may become activated by the mere observation of the first act of the chain,
273 the observer is able to predict future motor behaviour and the consequences
274 of the whole action sequence before its execution, i.e. the co-actor's motor
275 intention. However, since a single motor act may be part of several chains,
276 or may not be directly observable, the integration of additional contextual
277 information is necessary to disambiguate the simulation (Erlhagen et al.,
278 2007). For the assembly task, an important input comes from layer OML
279 representing the memorized world knowledge about the location of the
280 different parts in the two working areas. A second source of information
281 that may sustain the simulation process is the shared task knowledge about

282 what the human partner could do in a particular joint action situation
283 (Sebanz, Knoblich and Bekkering, 2006). The subgoals of the assembly
284 work, that are currently available for the team, are represented by
285 populations in the common subgoal layer (CSGL). They are continuously
286 updated in accordance with the assembly plan based on visual feedback
287 about the state of the construction and the inferred goal of the co-actor
288 (represented in the IL). The input from the IL to the CSGL is of particular
289 importance for pro-active behaviour since an updating of subgoals based on
290 anticipated action outcomes allows the robot to plan ahead of time a
291 complementary behaviour that best serves the user's future needs. The
292 CSGL contains two sublayers each containing populations representing all
293 possible subgoals of the assembly task. The sequential order of task
294 execution is encoded by the connections between populations in the two
295 layers. Input signalling the achievement (or predicted achievement) of a
296 certain subtask activates the respective population representation in the first
297 layer which in turn drives automatically through the connections the
298 populations in the second layer representing the next possible assembly
299 steps (e.g., attaching first a wheel and subsequently fixing it with a nut). The
300 action execution layer (AEL) contains the same goal-directed action
301 sequences as the ASL. The different populations integrate input from the IL,
302 OML and CSGL to select among all possible actions the most appropriate
303 complementary behaviour.

304 The implemented context-sensitive mapping from observed actions on to-be
305 executed complementary actions guarantees a fluent team performance if no
306 errors occur (Bekkering et al., 2009). To cope in an efficient manner also
307 with unexpected or erroneous behaviour of the co-actor, populations in the
308 error monitoring layer (EML) are sensitive to a mismatch on the goal level
309 (integrating input from CSGL and IL, e.g., the co-actor reaches a part that
310 has to be attached only later) or on the level of action means to achieve a
311 valid sub-goal (integrating input from OML and ASL, e.g., the co-actor
312 requests a certain part versus reaching the part directly in his/her
313 workspace). Through direct connections to the AEL, population activity in
314 the EML may bias the robot's planning and decision process by inhibiting
315 the representations of complementary actions normally linked to the inferred
316 goal and exciting the representations of a corrective response. Importantly,
317 to efficiently communicate detected errors to the human partner a corrective
318 response may consist of a manual gesture like pointing or a verbal comment
319 to attract the co-actor's attention (Bicho, Louro and Erlhagen, 2010).

320

321

322 **4. Basic concepts of the dynamic field framework**

323

324 The dynamics of each population in the various layers of the control
325 architecture is governed by a neural field equation (Wilson and Cowan,
326 1973; Amari, 1977). Dynamic neural field (DNF) architectures implement
327 the idea that task-relevant information is encoded by means of activation
328 patterns of local pools of neurons (Erlhagen and Bicho, 2006). These
329 patterns are initially triggered by input from connected populations and
330 sources external to the network (e.g., vision system, speech input). They
331 may become self-stabilized in the absence of any external input due to the
332 recurrent interactions within the population. Figure 3 shows an example of
333 a self-sustained activity pattern (dashed-dotted line) representing a grasping
334 behaviour. Importantly, there exists an instability of the field dynamics. The
335 self-stabilized pattern coexists with a stable homogenous activation
336 distribution (solid line) that represents the absence of specific information
337 about the motor primitive (resting level). Only sufficiently strong input may
338 activate the self-sustaining forces within the population. Weaker external
339 stimuli lead to a subthreshold, input-driven activation pattern (dashed line).
340 This preshaping of local populations by relative weak input signals may
341 nevertheless play an important role for the processing in the joint action
342 circuit. It brings populations closer to the threshold for triggering the self-
343 sustaining interactions and thus biases the decision processes linked to
344 behaviour (Erlhagen and Schöner, 2002).

345

346 -----Insert Figure 3 around here -----

347

348 We employed a particular form of a DNF first analyzed by Amari (1977). In
349 each model layer i the activity $u_i(x,t)$ at time t of a neuron at field location x
350 is described by the following integro-differential equation (for a discussion
351 of analytical results see Erlhagen and Bicho, 2006):

$$352 \quad \tau_i \frac{\delta u_i(x,t)}{\delta t} = -u_i(x,t) + S_i(x,t) + \int w_i(x-x') f_i(u(x',t)) dx' + h_i \quad (1)$$

353 where the constants $\tau_i > 0$ and $h_i < 0$ define the time scale and the resting
354 level of the field dynamics, respectively. The integral term describes the
355 intra-field interactions. It is assumed that (1) the interaction strength,
356 $w_i(x,x')$, between any two neurons x and x' depends only on the distance
357 between locations, and that (2) nearby cells excite each other, whereas
358 separated pairs of cells have a mutually inhibitory influence. For the present
359 implementation we used the following integral kernel of lateral-inhibition
360 type:

$$361 \quad w_i(\Delta x) = A_i \exp\left(-\Delta x^2 / (2\sigma_i^2)\right) - w_{inhib,i} \quad (2)$$

362 where $w_{inhib,i}>0$ is a constant and $A_i>0$ and $\sigma_i>0$ describe the amplitude and
 363 the standard deviation of a Gaussian, respectively. Only sufficiently
 364 activated neurons contribute to interaction. The threshold function $f_i(u_i)$ is
 365 chosen of sigmoidal shape with slope parameter β and threshold u_0 :

$$366 \quad f_i(u_i) = \frac{1}{1 + \exp(-\beta(u_i - u_0))}. \quad (3)$$

367 Normally, summed input from several connected populations is necessary to
 368 create a self-stabilized activity pattern in a target population. Each
 369 connected population contributes a Gaussian input signal of certain strength
 370 whenever its activity level is above threshold. The total input from all
 371 connected populations in the various layers of the dynamic field architecture
 372 layer u_i can thus be mathematically described by

$$373 \quad S_l(x, t) = K \sum_m \sum_j a_{mj} c_{ij}(t) \exp\left(-\frac{(x - x_m)^2}{2\sigma^2}\right) \quad (4)$$

374 where $c_{ij}(t)$ is a function that signals the existence or evolution of a self-
 375 sustained activation pattern in subpopulation j in layer u_l , and a_{mj} is the
 376 inter-field synaptic connection between subpopulation j in u_l to
 377 subpopulation m in u_i . The parameter K scales the total input to the target
 378 population relative to the threshold for triggering a self-sustained pattern.
 379 This guarantees that the inter-field coupling is weak and the field dynamics
 380 is dominated by the recurrent interactions.

381 The existence of a single, self-stabilized pattern of activation in a dynamic
 382 field can not only be used to implement a working memory function but is
 383 also closely linked to decision making. In layers ASL, IL, AEL and AML
 384 subpopulations encoding different action chains (ASL), goals (IL),
 385 complementary actions (AEL) and detected errors (EML), respectively,
 386 interact through lateral inhibition. These inhibitory interactions lead to the
 387 suppression of activity below resting level in competing neural pools
 388 whenever a certain subpopulation becomes activated above threshold.
 389 Figure 4 shows an example of the temporal evolution of activity in a field
 390 encoding different actions. The population for which the summed input
 391 from connected populations is highest wins the competition process. Note
 392 that at the beginning all subpopulations appear to be activated to some
 393 extent, that is, all action alternatives receive input from suprathreshold
 394 activity in connected pools. At time $t=0$ an additional input to the
 395 population encoding action A2 drives the activity beyond the critical level
 396 for a self-stabilized pattern.

397

398

-----Insert Figure 4 around here -----

399

400 To represent and memorize simultaneously (1) the location of several
401 objects, and (2) multiple common subgoals, the spatial ranges of the lateral
402 interactions in layers OML and CSGL were adapted to avoid a direct
403 competition between different populations. The updating of the memorized
404 information is performed by defining a proper dynamics for the inhibition
405 parameter, $h_i < 0$, of the population dynamics (Bicho, Mallet and Schöner,
406 2000). A sufficiently large global inhibition destabilizes an existing activity
407 peak. As a consequence, the population activity decays back to the stable
408 resting state.

409
410

411 **5. Results**

412

413 In the following we validate the dynamic field architecture by presenting
414 snapshots of the human-robot interactions in the assembly task. The
415 examples illustrate the impact of action observation on decision making in
416 varying context from the perspective of the robot. The focus is on showing
417 and explaining the goal inference, error detection and anticipatory action
418 selection capacities. In all examples here reported the sequential order of
419 sub-tasks for the construction of the ‘toy vehicle’ is the following: First,
420 mount wheels and fix them with nuts; second, insert column 1; third, insert
421 column 2 and column 4; fourth, insert column 3, and finally mount the top
422 floor.

423 The videos of the human-robot interaction and the associated dynamics of
424 the fields can be found at

425 <http://dei-s1.dei.uminho.pt/pessoas/estela/JASTvideosBookIS.htm>.

426 Alternatively, the videos may be seen at

427 <http://www.youtube.com/watch?v=A0qemfXnWiE> (video 1) and

428 <http://www.youtube.com/watch?v=7t5DLgH4DeQ> (video 2).

429

430

431 *5.1 Pro-active behavior and goal inference based on an anticipatory* 432 *model of action observation*

433

434 -----Insert Figure 5 around here -----

435 -----Insert Figure 6 around here -----

436 -----Insert Figure 7 around here -----

437 -----Insert Figure 8 around here -----

438 -----Insert Figure 9 around here -----

439

440 An important prerequisite for successful and fluent interaction is that both
441 team members must be committed to the fulfillment of the joint task. For the

442 robot this means that it should be able to initiate the assembly work and take
443 initiative whenever the human partner is taking too long to act. Moreover,
444 the robot should be able to show its commitment to the team by selecting an
445 action in anticipation of the consequences of the co-actor's action.
446 Anticipation of action effects is possible since the robot simulates the co-
447 actor's motor intentions using its own knowledge about goal-directed action
448 sequences in the assembly task. The two capacities are tested in the
449 experiment illustrated in Figure 5. The robot's camera view shows that all
450 components are distributed on the table (Fig. 6, panel B). The activity
451 patterns in the OML represent the knowledge about the corresponding
452 distribution of the components in the two working areas (Fig. 6, panels A
453 and C). The robot takes the initiative to start the assembly work while the
454 co-actor is still reading the instructions and thus does not show any object-
455 directed action. Since the robot has no wheel within its reach, it decides to
456 request a wheel to mount it on its side of the platform (Fig.5, snapshots S1-
457 S2). This decision is possible because the information about the available
458 subgoals (see the CSGL, Fig.7) and the location of parts in the two working
459 areas (see the OML, Fig.6) creates sufficient input to the AEL to trigger a
460 self-stabilized activation peak centered at the action 'request wheel' (see
461 Panel B in Fig.8, time interval T1-T2).

462 The user then grasps a wheel. However, her intention is not to transfer it to
463 the robot but to mount it on her side. As can be seen, at the moment of
464 grasping the wheel (Fig.5, snapshots S2-S3) the robot is able to anticipate
465 the partner's motor intention and immediately prepares for holding the base
466 in order to help the user while she inserts the wheel on the axle (Fig.5,
467 snapshots S4-S5). The capacity to infer the goal of the user at the time of
468 grasping is possible because of the way in which the partner grasps an
469 object conveys information about what she intends do with it. The robot has
470 sequences of motor primitives in its motor repertoire that associate the type
471 of grasping with specific final goals. A grasping from above is used to
472 attach a wheel to the axle whereas using a side grip is the most comfortable
473 and secure way to hand the wheel over to the co-actor. The observation of
474 an above grip (represented in the AOL) together with information about the
475 currently active subgoal(s) (attach wheel on the user's side) trigger an
476 activation peak in ASL that represents the simulation of the corresponding
477 'reaching-grasping-inserting' chain (see Panel A in Fig.8, time interval T2-
478 T3), which automatically activates the underlying goal in the intention layer
479 (see Fig.9, T2-T3). The evolving activation pattern in the AEL (panel B,
480 Fig.8, T2-T3) reflects the decision to stabilize the base. The robot's
481 decision to give up its own intention to attach a wheel is a result of slight
482 differences in the connection strengths between representations in the CSGL

483 and the AEL. These differences favor the realization of the user’s subtasks
484 over the subtasks that are under the control of the robot.

485 When the user has attached the wheel the corresponding activation peak in
486 CSGL disappears and an activation pattern representing the subsequent
487 subgoal (‘insert nut on user’s side’ to fix the wheel) automatically evolves
488 (see Fig.7, T3-T4). The second subgoal ‘insert wheel on robot’s side’
489 remains active.

490 The user again takes long to act, this time because she is interrupted by a
491 colleague entering the room (see video 1). The robot again takes the
492 initiative and demands a wheel (Fig.5, snapshots S6-S7) since as before the
493 information represented in the CSGL and the OML is sufficient to trigger
494 the corresponding action representation in the AEL (Panel B in Fig.8, T3-
495 T4). Next, the user grasps a wheel with a ‘side grip’ and the robot
496 anticipates that she is going to hand it over (see ASL in Panel A in Fig.8,
497 T4-T5). The robot prepares to receive the wheel in order to mount it on its
498 side of the platform (see snapshots S7-S9 in Fig.5, and the AEL activation
499 in Fig.8, T4-T5).

500

501

502 *5.2 Understanding partially occluded actions and anticipating the user’s* 503 *future needs*

504

505 -----Insert Figure 10 around here -----

506 -----Insert Figure 11 around here -----

507 -----Insert Figure 12 around here -----

508 -----Insert Figure 13 around here -----

509

510 In the previous example we have seen that the robot could infer through
511 motor simulation the co-actor’s motor intention from the way the object is
512 grasped. But what happens when the robot cannot directly observe the
513 hand-object interaction? In natural environments with multiple objects and
514 occluding surfaces this is a realistic scenario. The capacity to discern the
515 user’s motor intention and to select an appropriate complementary behavior
516 should of course not be disrupted by missing information about the grip type
517 used. Information about the context in which the action is executed may
518 sustain the motor simulation process. This is illustrated in the following
519 interaction scenario in which only the ‘reaching’ part of the user’s action
520 sequence can be observed. The robot sees the hand disappearing behind an
521 occluding surface but knows that there is a wheel behind (Fig.10) since the
522 occluder has been introduced into the scene only after the robot could
523 memorize the position of the wheel in the workspace. Figure 11 illustrates
524 the goal inference mechanism in this situation. The AOL (not shown) only

525 codes the reaching behavior. The currently possible subgoals represented in
526 CSGL are ‘insert wheel on user’s side’ and ‘insert wheel on robot’s side’
527 (Fig.11, Panel A). The inputs from the AOL and the CSGL to the ASL are
528 thus compatible with two competing chain representations and thus should
529 only pre-activate these representations. The additional input necessary for
530 goal inference comes from the information about the location of the wheel
531 in the user’s workspace represented in the OML (see Panel B in Fig. 11).
532 This input triggers the evolution of a self-stabilized activation peak in the
533 ASL representing the action sequence ‘reach wheel-grasp-insert’ (see Panel
534 C in Fig.11). This activation in turn induces a suprathreshold pattern in the
535 IL representing the underlying goal ‘insert wheel’ (see Panel B in Fig.12;
536 see also snapshot S4 in Fig.10). When the activation pattern in IL rises
537 above threshold it initiates a dynamic updating process in the second layer
538 of the CSGL, representing the next possible subgoal(s) for the user (Panel
539 C, Fig.12; see also snapshot S5 in Fig. 10). This representation allows the
540 robot to select a complementary action that serves user’s future needs.
541 In summary, in this example ‘insert wheel on robot side’ and ‘insert wheel
542 on user’s side’ are two currently available subgoals, the robot infers that the
543 user is grasping a wheel with the intention to mount it. The next possible
544 subgoal for the user is ‘insert nut’ in order to fix the wheel. All the nuts and
545 one wheel are located in the robot’s workspace. As a consequence, three
546 complementary actions in AEL are supported by input from connected
547 populations and thus compete for expression in the robot’s overt behavior
548 (see Fig.13): ‘insert wheel’, ‘give nut’ and ‘hold the base’. As can be seen
549 when comparing the pattern of activation that evolves in the AEL, the robot
550 decides to serve the human by grasping a nut for handing it over (see Fig.10,
551 snapshots S5-S6).

552
553

554 *5.3 Error detection and context-sensitive interpretation of a request* 555 *gesture*

556
557
558
559
560

-----Insert Figure 14 around here -----
-----Insert Figure 15 around here -----
-----Insert Figure 16 around here -----

561 Even in known tasks the user can easily become confused and make errors
562 that should be corrected by the co-actor before failure becomes manifested.
563 For example, the user may get confused with the different columns (labeled
564 as C1 to C4) and tries to insert a certain column in the wrong whole or tries
565 to manipulate the columns in a wrong sequential order. The example shown
566 in Fig.14, illustrates two cases in which the robot’s error monitoring layer

567 detects mismatches between the inferred intention of the user and the state
568 of the construction, that is, between the intention and possible subgoals. In
569 the two cases, the robot reports the error to the user and explains what needs
570 to be done. This interaction scenario also allows illustrating the robot's
571 ability to infer the goal of a request gesture.

572 As illustrated by snapshot S16 (in Fig.14), the robot observes the user
573 grasping column C4 with a top grip which it interprets via action simulation
574 as belonging to a 'grasping to insert sequence' (see Fig.15, time interval T7-
575 T8). However, the activation pattern in the DNF of the CSGL representing
576 present available subgoal(s) indicates a conflict. The correct subtask is to
577 'insert column C1 on the robot's side' (Panel A in Fig.16, time interval T7-
578 T8). Input from this field together with input from IL (Panel B in Fig.16,
579 time interval T7-T8) triggers a suprathreshold peak of activation in a
580 population encoding the error in user's intention (see Panel C in Fig.16,
581 time interval T7-T8). The error related activity is linked to a subpopulation
582 in AEL that represents a corrective response. In this case, suprathreshold
583 activity initiate speech output to report the mismatch and to explain what
584 needs to be done (see snapshot S17-S18 in Fig.14, and field activity in the
585 AEL depicted in Panel D of Fig.16). The content of the speech combines the
586 information coded in the activation patterns that have initially triggered the
587 error related activity, i.e. intention layer (IL) and possible active sub-goals
588 represented in CSGL, respectively ("*You cannot insert column 4 yet. First,*
589 *we need to insert column 1*"). Subsequently, the user proceeds by grasping
590 column C1 (snapshot S19, Fig.14) which the robot interprets as belonging to
591 a 'grasp to handover' chain (see Fig. 15, time interval T8-T9). An activation
592 pattern starts to evolve in AEL (Panel D in Fig. 16, time interval T8-T9) that
593 represents the robot's decision to receive column 1 for inserting it in the
594 corresponding whole of the platform (snapshots S20-S21, Fig.14).

595 Next the user opens up her empty hand as it moves towards the robot
596 (snapshot 22, Fig.14). The robot has this gesture associated with 'request
597 object' in its motor repertoire. The observation of this unspecific gesture
598 activates to some extent all action chains in the ASL linked to requesting the
599 components of the toy vehicle that are in the robot's workspace (in this
600 example nut and column C3). The nuts have already been attached (working
601 memory about already achieved subtask is represented by self-stabilized
602 activation patterns in the "past" layer of the CSGL, not shown). Thus the
603 robot interprets the user's gesture as a request for column C3 (see snapshot
604 S22 in Fig14, and field activity in ASL in Fig.15, time interval T9-T10).
605 The bimodal activation pattern in CSGL (panel A, Fig.16) indicates that the
606 currently available subgoals are inserting columns C2 and C4 on the user's
607 side of the platform. Hence similar to the preceding example, input from
608 this field together with input from IL induces a suprathreshold peak of

609 activation in a population of the error layer representing that ‘inserting
610 column C3’ is an error in intention (see snapshot S23 in Fig.14, and panel D
611 in Fig.16, in the interval T9-T10). The robot verbally signals the error and
612 explains what sub-goals are currently possible (snapshot S24, Fig.14).
613 Subsequently, the user grasps and inserts column C4 while the robot assists
614 the human user by stabilizing the base.

615

616

617 **6. Discussion**

618

619 We presented a robot control architecture for human-robot collaboration that
620 is inspired by neurocognitive theories about how humans perceive and act in
621 a social context. The results of the validation in a joint assembly task show
622 that the implementation of a human-like joint action model in the robot
623 supports a fluent and flexible task execution. The ease with which humans
624 coordinate in routine joint activity their decisions in space and time is
625 impressive. The capacity to quickly register the co-actor's motor intention
626 before his or her action sequence is completed is essential for a fluent
627 performance in human team activity (Sebanz, Knoblich and Bekkering,
628 2006). Being able to select an action based on predicted effects of the co-
629 actor's behaviour is thus considered a crucial skill for robots in order to be
630 fully accepted by a human user as a social partner in cooperative tasks
631 (Hoffman and Breazeal , 2007). Converging lines of experimental evidence
632 support the notion of an automatic and obligatory motor simulation process
633 as the underlying mechanism for the intention understanding capacity (for
634 review see Rizzolatti and Craighero, 2004).

635 The dynamic field architecture implements the idea that goal inference,
636 performance monitoring and action selection occur rather effortlessly and do
637 not require a fully developed human capacity for conscious control (for a
638 review Ferguson and Bargh, 2004). As the representation of context, goals
639 and shared task knowledge are interconnected, the observation of a motor
640 act together with situational cues activates through motor simulation first
641 the self-sustained population representations of the related goal and
642 subsequently the representation of the most appropriate complementary
643 action. As our examples show, this automatic process includes situations in
644 which the human partner acts in an unexpected or inappropriate manner.
645 This view on intention communication and joint action planning contrasts
646 with most robot control architectures that have been tested in the past in
647 similar collaborative tasks (e.g., Alami et al., 2005; Steil et al., 2004; Gast et
648 al., 2009; but see Breazeal, Gray and Berlin, 2009, for a conceptually
649 similar approach). Typically, these architectures include dedicated modules
650 that organize the high-level task of intention coordination between the co-

651 actors and the intelligent monitoring of the team performance using some
652 form of symbolic manipulation and logic. Although we do not deny that the
653 results of our robotics experiments could be implemented in a symbolic
654 framework we argue here that the additional planning process that would be
655 needed to link the high-level representations to the motor level of the robot's
656 actuators would greatly reduce the effectiveness of those representations.
657 In the dynamic field architecture the decision process linked to
658 complementary actions unfolds over time under multiple influences which
659 are themselves modelled as dynamic representations with proper time
660 scales. This is the basis of flexible behaviour in dynamic joint action
661 conditions. The absence or delay of information about for instance the co-
662 actor's motor intention will automatically lead to a decision that does not
663 take into account the co-actor (Bicho et al., in press). Conversely, the
664 dynamic updating of the currently available subgoals for the team based on
665 predicted effects of the co-actor's ongoing action allows for anticipatory
666 action planning. If on the other hand the predicted effect is inconsistent with
667 the current goals for the team the obligatory integration of evolving
668 suprathreshold activity in the action monitoring layer may override the
669 planning of a pre-potent complementary action.

670 Among the many possible types of errors that might occur during joint
671 activity of the human-robot team (for a discussion see Spexard et al., 2008)
672 we have focused in this chapter on the detection and communication of
673 intention errors made by the human partner. The proposed action monitoring
674 mechanism can be easily extended to cope with other types of unexpected
675 events including errors made by the robot. For instance, the co-actor may
676 show a request gesture with the intention to fulfil the valid subgoal of fixing
677 the wheel with a nut. However, since a nut is located in her workspace,
678 responding to the request by transferring the nut would not be an efficient
679 behaviour of the robot. Population activity in the EML that combines the
680 information from the OML about object location and the ASL about the
681 goal of the request gesture should instead trigger for instance a pointing
682 gesture to attract the co-actor's attention to the nut in her workspace.
683 Similarly, population activity in the AML that automatically integrates the
684 information about the goal (represented in CSGL) of a self-performed
685 sequence like reaching-grasping-attaching a wheel and proprioceptive
686 (and/or visual) information about an accidental loss of the wheel during
687 transportation may allow the robot to quickly start searching for a new
688 wheel or to ask the human for assistance. Interestingly, studies investigating
689 the functional system for action monitoring suggest that humans use similar
690 cognitive and neural mechanisms to detect own and observed errors in joint
691 action (Bekkering et al., 2009).

692 The focus of the present experiments was on implicit communication. The
693 robot had to interpret the co-actors actions and gestures in terms of a goal to
694 select an adequate complementary action. Although we have argued that
695 joint activity in a familiar task often does not require explicit
696 communication it is undeniable that being able to communicate with the
697 robot system by natural speech would greatly facilitate human-robot
698 interactions in many situations (e.g., Spexard et al., 2007; Pardowitz et al.,
699 2007; Gast et al., 2009). Our robot speaks aloud to make its goal inference
700 and action monitoring capacities transparent for the user. The knowledge
701 about the robot's cognitive capacities supports predictability of its behaviour
702 which is essential to an effective collaboration. Using a simplified version
703 of a joint assembly task in which the robot merely assists the human user by
704 handing over pieces (Bicho, Louro and Erlhagen, 2010), we have recently
705 made first steps towards integrating in the DNF-model of joint action the
706 capacity to understand simple action-related speech. The basic idea is that
707 the automatic resonance of motor structures during action observation
708 extends to the language domain. The robot understands sentences like *Give*
709 *me the wheel* or *I give you the wheel* by covertly activating semantically
710 congruent motor representations that are linked to the specific goal or end-
711 state (e.g. hand opens up as it moves towards the co-actor for a request
712 gesture and a reaching-grasping-holding out sequence for a handing over
713 procedure). This embodied view on language comprehension is supported
714 by findings in a range of recent experimental studies (for review see Fischer
715 and Zwaan, 2008).

716
717

718 **7. Conclusions and outlook**

719

720 The work presented in this chapter wished to contribute to the development
721 of design principles for robots that are supposed to directly collaborate with
722 their human partners in shared tasks. As an exquisitely social species,
723 humans are experts in coordinating actions and decisions with other in order
724 to achieve common goals. We believe that implementing a human-like joint
725 action model in the robot is a promising approach because it will allow the
726 artificial cognitive agent to meet the user's expectations about a pleasant,
727 efficient and successful interaction with a socially intelligent partner. While
728 theories about the neurocognitive mechanisms supporting human joint
729 action have now reached a sufficient level of detail to guide robotics work,
730 it is also clear that these theories contain hypothesis and assumptions for
731 which the experimental evidence is still under debate. Implementing and
732 testing theories and hypothesis about human joint action in an embodied
733 agent with sensory, motor and cognitive capabilities offers in our view

734 unique opportunities for researchers from cognitive science and
735 neuroscience.

736 The adopted dynamic perspective offers in general a high degree of
737 flexibility and robustness in joint task execution. However, the present
738 implementation of the dynamic field architecture limits the cooperative
739 interactions to the specific assembly task since the neural representations
740 and their connectivity were tailored by the designer. It is thus highly
741 desirable to endow the robot with a developmental program that would
742 allow the artificial agent to autonomously learn and represent new task-
743 relevant representations (Weng, 2004). Learning efficient joint action
744 coordination in a complex task is a very demanding and to a large extent
745 unsolved problem even when starting with a "minimal" set of pre-defined
746 capacities and knowledge. Adopting a socially guided machine learning
747 paradigm in which a human trainer teaches a robot through demonstration
748 and verbal or gestural commands in much the same way as parents teach
749 their children seems to be a promising research direction (Otero et al., 2008;
750 Thomaz and Breazeal, 2008). First experimental results of our attempt to
751 apply a learning dynamics for establishing inter-field connections show the
752 feasibility of the approach. Using correlation-based learning rules (Gerstner
753 and Kistler, 2002) with a gating that signals the success of behaviour, we
754 have shown for instance how goal-directed mappings between action
755 observation and action execution that support an action understanding
756 capacity may develop during learning and practice (Erlhagen, Mukovskiy
757 and Bicho, 2006; Erlhagen et al., 2006). Importantly, the developmental
758 process may explain the emergence of new task-specific populations which
759 have not been introduced to the architecture by the human designer
760 (Erlhagen et al., 2007).

761 We are currently applying and testing a learning by demonstration approach
762 to systematically address the question of how ARoS may acquire and store
763 the knowledge about the serial order of task execution that was predefined
764 by the designer in the experiments reported here. For the present dynamic
765 field architecture this means to autonomously develop the connections
766 between neural populations in the two sublayers of CSGL representing
767 subsequent steps of the assembly plan. We exploit here the self-stabilizing
768 properties of the field dynamics determined by the recurrent interactions
769 within each population and fixed excitatory and inhibitory connections
770 between neuronal populations in both layers representing the same assembly
771 step or subgoal of the construction plan. During demonstration by a human
772 teacher, ARoS perceives changes in the state of the construction. This visual
773 input triggers a suprathreshold activity pattern of the neural population in
774 the second layer representing the currently achieved subgoal. Mediated by
775 excitatory connections this activity propagates to the population in the first

776 layer and initiates the evolution of a self-sustained activity pattern. This
777 pattern implements a working memory function in the first layer and
778 suppresses through inhibitory feedback connections the population activity
779 in the second layer. Importantly, the actively maintained representation of
780 the achieved subgoal allows the robot to learn associations between
781 subsequent assembly steps that are separated in time. Correlation-based
782 learning takes place whenever a perceived change in the state of
783 construction triggers a transient population representation of a newly
784 achieved subgoal in the second layer. Since the serial order of task
785 execution may not be exactly the same in different demonstrations of the
786 task (e.g., different teachers), association between a single memorized
787 subgoal and several possible next steps may be learned during observation.
788 The work on learning and development in the dynamic field architecture for
789 joint action represents first steps towards robotics systems that will
790 ultimately be able to autonomously built representations for assisting
791 different human users in a large variety of tasks. We believe that combining
792 the processing principles of neural field dynamics and different machine
793 learning techniques in the context of the socially guided learning paradigm
794 represents a promising research direction towards achieving this demanding
795 goal.

796

797 **Acknowledgements**

798 The present research was conducted in the context of the fp6-IST2 EU-IP
799 Project JAST (proj. nr. 003747) and partly financed by the FCT grants
800 POCI/V.5/A0119/2005 and CONC-REEQ/17/2001. We would like to thank
801 Emanuel Sousa, Flora Ferreira and Toni Machado for their help during the
802 robotic experiments.

803

804 **References**

805

806 Alami, R., Clodic, A., Montreuil, V., Sisbot, E. A. & Chatila, R. (2005).
807 Task planning for human-robot interaction. In: *Proceedings of the 2005*
808 *Joint Conference on Smart Objects and Ambient Intelligence*, ACM
809 International Conference Proceeding Series, 121, 81-85.

810

811 Allisandrakis, A., Nehaniv, C. L., & Dauthenhahn, K. (2002). Imitation with
812 ALICE: learning to imitate corresponding actions across dissimilar
813 embodiments. *IEEE Transactions on Systems, Man and Cybernetics-Part A*,
814 Syst. Humans, 32, 482-492.

815

816 Amari, S. (1977). Dynamics of pattern formation in lateral-inhibitory type
817 neural fields. *Biological Cybernetics*, 27, 77–87.

818

819 Bekkering, H., Bruijn, E.R.A. de, Cuijpers, R.H., Newman-Norlund, R.,
820 Schie, H.T. van & Meulenbroek, R.G.J. (2009). Joint action: neurocognitive
821 and neural mechanisms supporting human interaction. *Topics in Cognitive*
822 *Science*, 1, 340-352.

823

824 Bicho, E., Mallet, P. & Schöner, G. (2000). Target representation on an
825 autonomous vehicle with low-level sensors. *International Journal of*
826 *Robotics Research*, 19, 424-447.

827

828 Bicho, E., Louro, L., Hipolito, N. & Erlhagen, W. (2009). A dynamic field
829 approach to goal inference and error monitoring for human-robot
830 interaction. In: Dautenhahn, K. (Ed.), Proceedings of the 2009 *International*
831 *Symposium on New Frontiers in Human-Robot Interaction*. AISB 2009
832 Convention, Heriot-Watt University Edinburgh, 31–37.

833

834 Bicho, E., Louro, L. & Erlhagen, W. (2010). Integrating verbal and
835 nonverbal communication in a dynamic neural field architecture for human-
836 robot interaction. *Frontiers in Neurorobotics*, 4, May, 2010. doi:
837 10.3389/fnbot.2010.0005.

838

839 Bicho, E., Louro, L., Costa e Silva, E. & Erlhagen, W. (in press). Neuro-
840 cognitive mechanisms of decision making in joint action: a Human-robot
841 interaction study. *Human Movement Science*.
842 doi:10.1016/j.humov.2010.08.012

843

844 Billard, A, Calinon, S., Dillmann, R. and Schaal, S. (2008). Robot
845 Programming by demonstration. In Handbook of Robotics. B. Siciliano, B.
846 and Kathib, O. (Eds). Springer.

847

848 Breazeal, C. (2004). Social interactions in HRI: the robot view. *IEEE*
849 *Transactions on Systems, Man and Cybernetics-Part C: Applications and*
850 *Reviews*, 34, 181-186.

851

852 Breazeal, C., Gray, J. & Berlin, M. (2009). An embodied cognition
853 approach to mindreading skills for socially intelligent robots. *International*
854 *Journal of Robotics Research*, 28, 656–680.

855

856 Costa e Silva, E., Bicho, E., Erlhagen, W & Meulenbroek, R. (submitted).
857 Human-like collision-free arm movements for human-robot collaboration.

858
859 Demiris, Y., Johnson, M. (2003). Distributed, predictive perception of
860 actions: a biologically inspired robotics architecture for imitation and
861 learning. *Connection Science*, 15, 231-243.
862
863 Erlhagen, W. & Schöner, G. (2002). Dynamic field theory of movement
864 preparation. *Psychological Review*, 109, 545-572.
865
866 Erlhagen, W. & Bicho, E. (2006). The dynamic neural field approach to
867 cognitive robotics. *Journal of Neural Engineering*, 3, R36–R54.
868
869 Erlhagen, W., Mukovskiy, A. & Bicho, E. (2006). A dynamic model for
870 action understanding and goal-directed imitation. *Brain Research*, 1083,
871 174–188.
872
873 Erlhagen, W., Mukovskiy, A., Bicho, E., Panin, G., Kiss, C., Knoll, A., van
874 Schie, H. & Bekkering, H. (2006). Goal-directed imitation for robots: a bio-
875 inspired approach to action understanding and skill learning. *Robotics and*
876 *Autonomous Systems*, 54, 353–360.
877
878 Erlhagen, W., Mukovskiy, A., Chersi, F. & Bicho, E. (2007). On the
879 development of intention understanding for joint action tasks. In:
880 *Proceedings of the 6th IEEE Int. Conf. on Development and Learning*.
881 Imperial College London, pp. 140–145.
882
883 Ferguson, M. J. & Bargh, J. A. (2004). How social perception can
884 automatically influence behaviour. *Trends in Cognitive Sciences*, 8, 33-39
885
886 Fischer, M. & Zwaan, R. (2008). Embodied Language: A Review of the
887 Role of the Motor System in Language Comprehension. *The Quarterly*
888 *Journal of Experimental Psychology*, 61, 825-850.
889
890 Fogassi, L., Ferrari, P. F., Gesierich, B., Rozzi, S., Chesri, F. & Rizzolatti,
891 G. (2005). Parietal lobe: from action organization to intention
892 understanding. *Science*, 308, 662-667.
893
894 Fong, T., Nourbakhsh, I. & Dautenhahn, K. (2003). A survey of socially
895 interactive robots. *Robotics and Autonomous Systems*, 42, 143–166.
896
897 Gast, J., Bannat, A., Rehrl, T., Wallhoff, F., Rigoll, G., Wendt, C., Schmidt,
898 S., Popp, M. & Farber, B. (2009). Real-time framework for multimodal
899 human–robot interaction. In: *Proceedings of the 2nd Conference on Human*

900 *System Interactions 2009, HIS '09. IEEE Computer Society, Catania, Italy,*
901 276–283.

902

903 Gerstner, W. & Kistler, W. M. (2002). Mathematical formulations of
904 Hebbian learning. *Biological Cybernetics*, 87, 404-415.

905

906 Hoffman, G. & Breazeal, C. (2007). Cost-based anticipatory action selection
907 for human-robot fluency. *IEEE Transactions on Robotics*, 23, 952–961.

908

909 Otero, N. Saunders, J., Dautenhahn, K. & Nehaniv, C. L. (2008). Teaching
910 Robot Companions: The Role of Scaffolding and Event Structuring.
911 *Connection Science*, 20, 111-134.

912

913 Pardowitz, M., Knopp, S., Dillmann, R. & Zöllner, R. D. (2007).
914 Incremental learning of tasks from user demonstrations, past experiences
915 and vocal comments. *IEEE Transactions on Systems, Man, and Cybernetics*
916 - *Part B: Cybernetics*, 37, 322-332.

917

918 Poljac, E., van Schie H.T. & Bekkering, H. (2009). Understanding the
919 flexibility of action-perception coupling. *Psychological Research*, 73, 578-
920 586.

921

922 Rizzolatti, G. & Craighero, L. (2004). The mirror-neuron system. *Annual*
923 *Review of Neuroscience*, 27, 169–192.

924

925 Schaal, A. (1999). Is imitation learning the route to humanoid robots?.
926 *Trends in Cognitive Science*, 3, 233-242.

927

928 Schöner, G. (2008). Dynamical systems approaches to cognition. In: Sun,
929 R. (Ed.). *The Cambridge Handbook of Computational Psychology*.
930 Cambridge University Press, pp. 101–125.

931

932 Sebanz, N., Bekkering, H. & Knoblich, G. (2006). Joint action: bodies and
933 minds moving together. *Trends in Cognitive Sciences*, 10, 70–76.

934

935 Silva, R., Bicho, E. & Erlhagen, W. (2008). ARoS: an anthropomorphic
936 robot for human-robot interaction and coordination studies. In: *Proceedings*
937 *of the 8th Portuguese Conference on Automatic Control (Control 2008)*. 21-
938 23 July, Vila Real, Portugal, 819–826.

939

940 Spexard, T. P., Hanheide, M., & Sagerer, G. (2007). Human-oriented
941 interaction with an anthropomorphic robot. *IEEE Transactions on Robotics*,
942 23, 852-862.

943

944 Spexard, T.P., Hanheide, M., Li, S. & Wrede, B. (2008). Oops, something is
945 wrong: Error detection and recovery for advanced human-robot interactions.
946 In: *Proceedings of the ICRA 2008 Workshop on Social Interaction with*
947 *Intelligent Indoor Robots*. Kruijff, G.J., Zender, H. Handeide, M. & Wrede,
948 B. (Eds).

949

950 Steil, J. J., Röthling, F., Haschke, R., & Ritter, H. (2004). Situated robot
951 learning for multi-modal instruction and imitation of grasping. *Robotics and*
952 *Autonomous Systems*, 47, 129-141.

953

954 Thomaz, A. L., & Breazeal, C. (2008). Teachable robots: Understanding
955 human teaching behaviour to built more effective robot learners. *Artificial*
956 *Intelligence*, 172, 716-737.

957

958 van Schie, H. T., van Waterschoot, B. M. & Bekkering, H. (2008).
959 Understanding action beyond imitation: Reversed compatibility effects of
960 action observation in imitation and joint action. *Journal of Experimental*
961 *Psychology: Human Perception and Performance*, 34, 1493–1500.

962

963 Weng, J. (2004). Developmental robotics: Theory and experiments.
964 *International Journal of Humanoid Robotics*, 1, 199-236.

965

966 Westphal, G., von der Malsburg, C. & Würtz, R. P. (2008). Feature-driven
967 emergence of model graphs for object recognition and categorization. In:
968 Kandel, A., Bunke, H., Last, M. (Eds.). *Applied Pattern Recognition*.
969 Springer Verlag, pp. 155–199.

970

971 Wilson, H. R., Cowan, J. D. (1973). A mathematical theory of the functional
972 dynamics of cortical and thalamic nervous tissue. *Kybernetik*, 13, 55–80.

973

974 Wilson, M. & Knoblich, G. (2005). The case for motor involvement in
975 perceiving co-specifics. *Psychological Bulletin*, 131, 460–473.

976

977

978

979

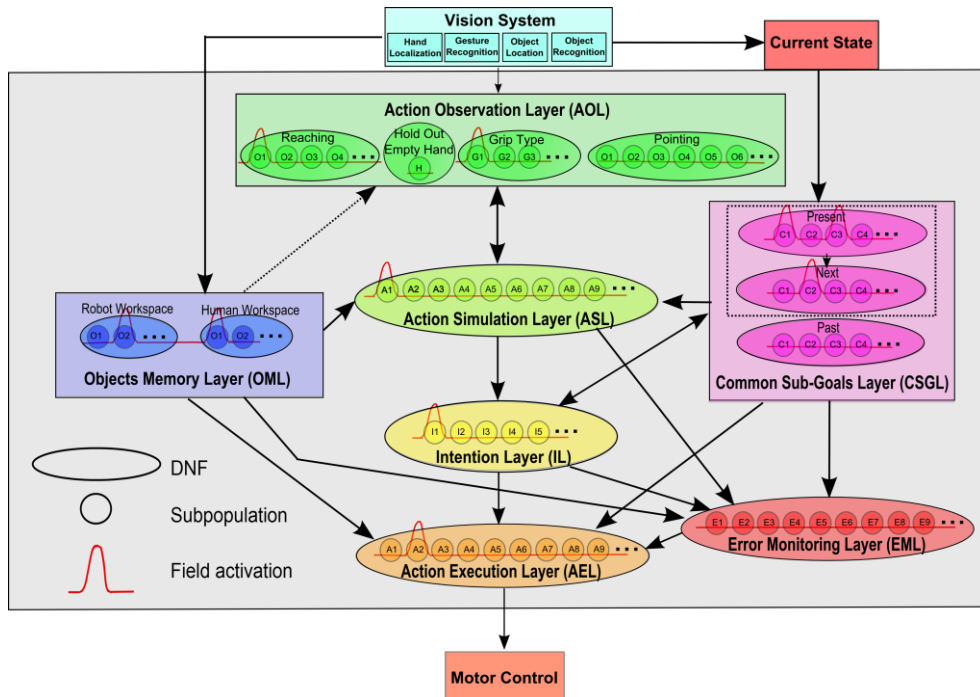
List of Figures & captions

980



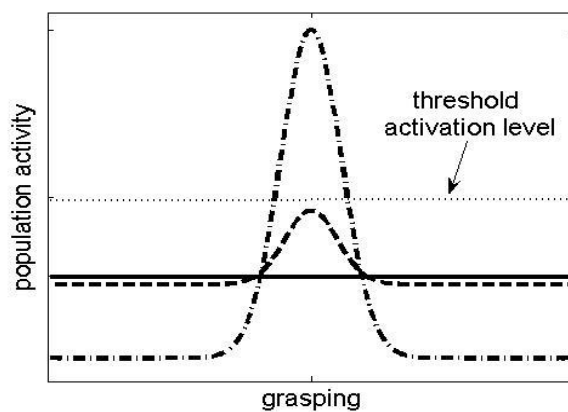
981

982 **Figure 1:** Human-robot team and scenario for the joint construction task.
983 The team has to collaborate in the construction of a 'toy vehicle' (shown in
984 panel b) from components that are initially distributed on a table (panel a).
985



986
 987
 988
 989
 990
 991
 992
 993
 994
 995

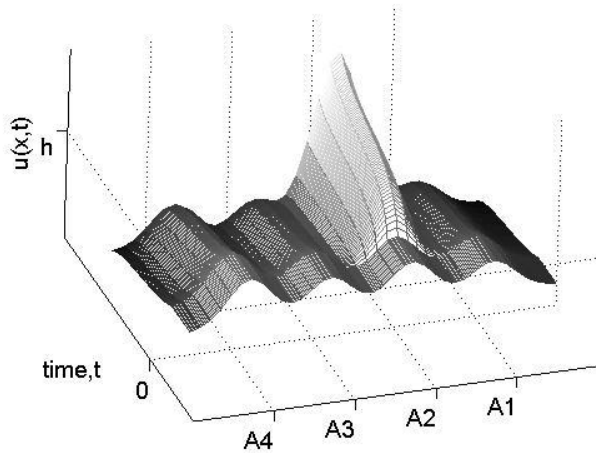
Figure 2: Schematic view of the cognitive architecture for joint action. It implements a flexible mapping from observed actions (layer AOL) onto complementary actions (layer AEL) taking into account the inferred action goal of partner (layer IL), detected errors (layer EML), contextual cues (OML) and shared task knowledge (CSGL). The goal inference capacity is based on motor simulation (layer ASL).



996
 997
 998
 999
 1000

Figure 3: The activity of a neural population is shown that represents through a self-stabilized activation peak the presence of information about a grasping behaviour (dashed-dotted line), while a flat, low-level distribution

1001 of activity (solid line) indicates that information about this motor primitive
1002 is currently not processed. In response to a weak input the population
1003 generates an activation pattern with amplitude below the threshold
1004 necessary to trigger an interaction-dominated activation peak (dashed line).
1005



1006 **Figure 4:** Decision making in a field representing different actions, A1 to
1007 A4. The decision is triggered by a sufficiently strong input at time $t=0$ to
1008 population A2. Note that at this time all 4 populations are activated below
1009 the critical value for a self-stabilized pattern.
1010

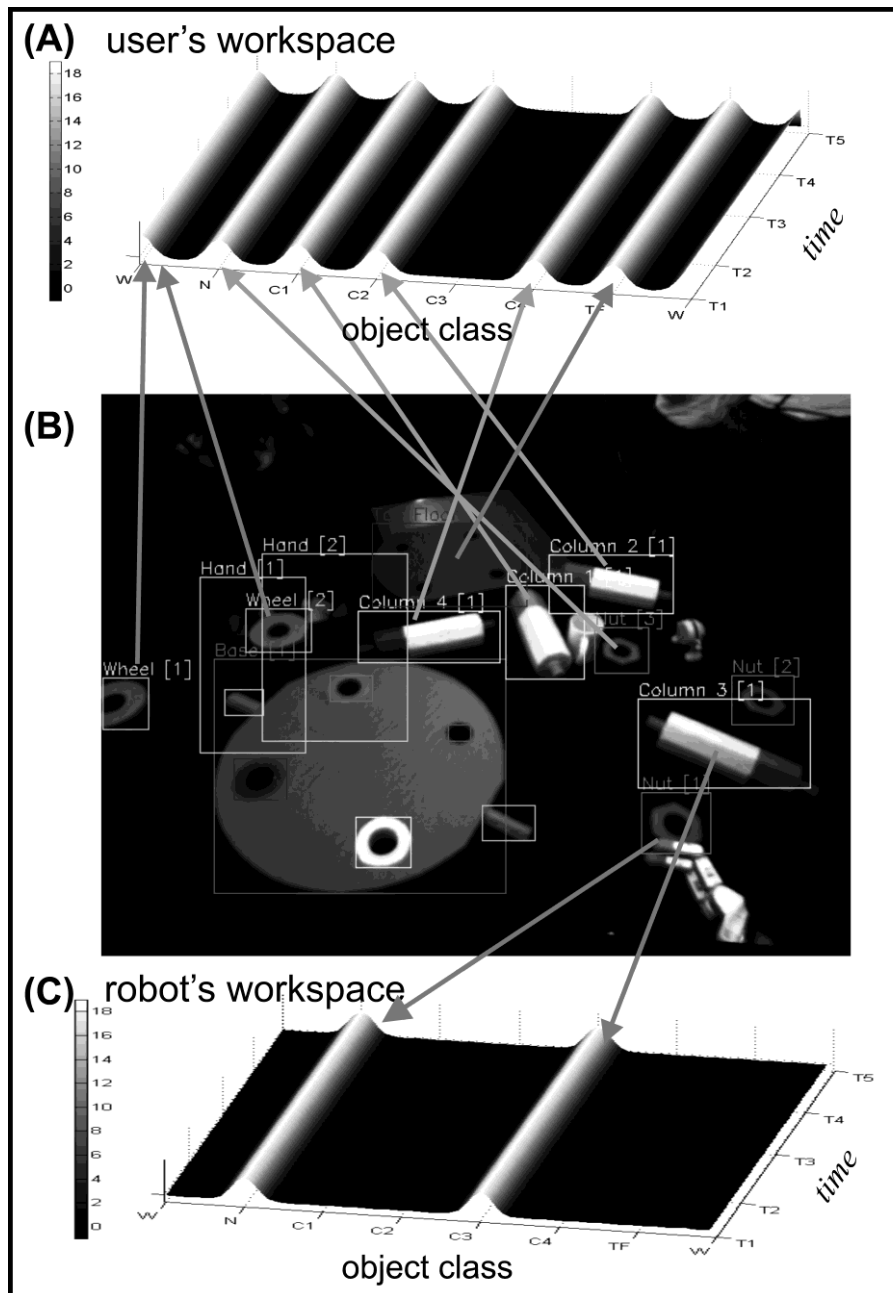
1011

1012



1013

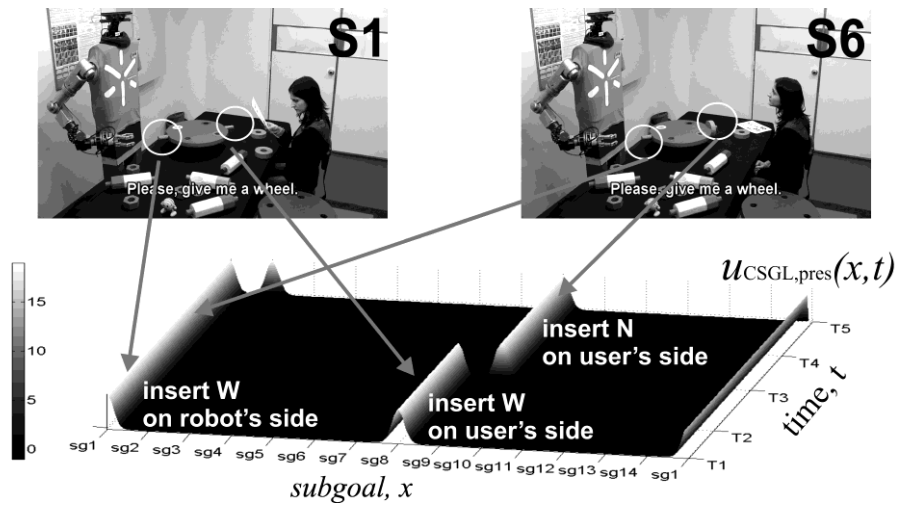
1014 **Figure 5:** Snapshots, S1-S9, in the time interval T1-T5 of video 1 (long
 1015 interaction scenario) illustrate the robot's pro-active behavior and its goal
 1016 inference capacity which is based on an anticipatory model of action
 1017 observation.
 1018



1019

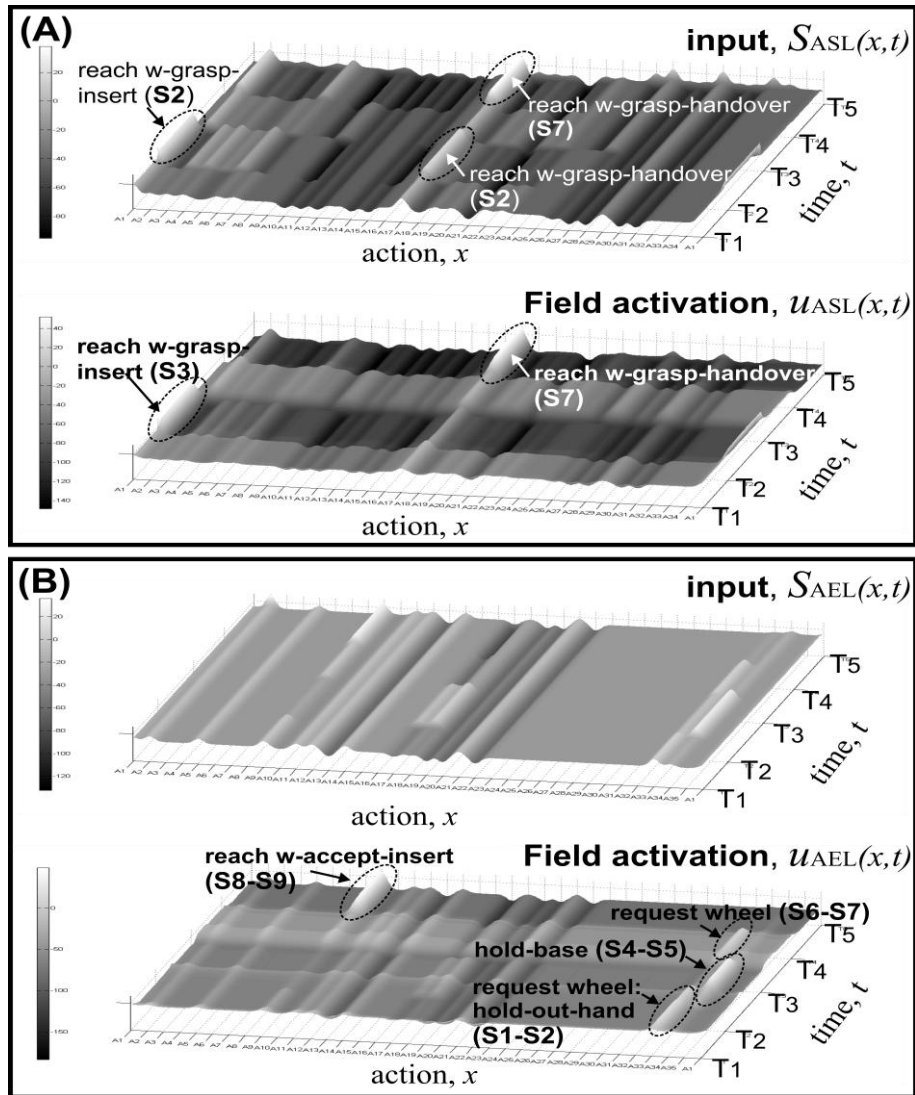
1020 **Figure 6:** Snapshot of the robot's vision system (Panel B) and
 1021 corresponding representations of objects in the OML. It contains two fields,
 1022 one for each workspace (i.e. Panels A and C). In each field the presence of
 1023 an object of a particular class is represented by a peak of activation.

1024



1025

1026 **Figure 7:** The temporal evolution of activity in the dynamic field encoding
1027 the currently available subgoals (in the CSGL) is shown in the interval T1-
1028 T5 (video 1). The snapshots S1 and S6 show the corresponding events of the
1029 human-robot interactions.



1030

1031

1032 **Figure 8:** Proactive behavior and goal inference based on an anticipatory

1033 model of action observation (video 1, time interval T1-T5). (A) Temporal

1034 evolutions of input to ASL (top) and field activity in ASL (bottom). (B)

1035 Temporal evolutions of input to AEL (top) and field activity in AEL

1036 (bottom)

1037

1038

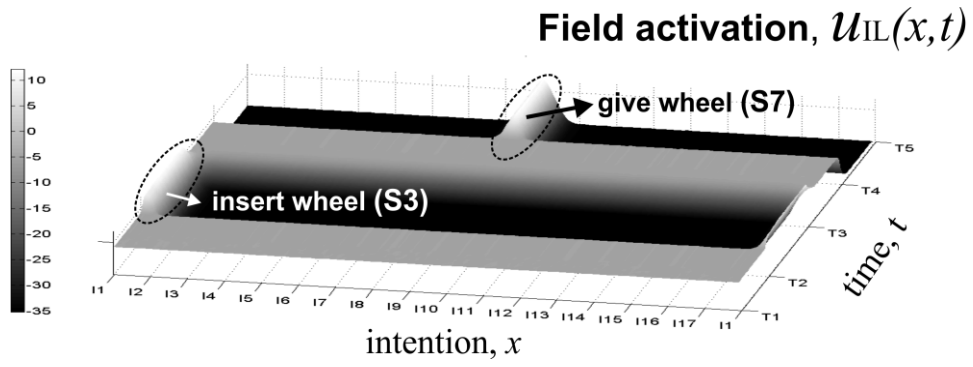
1039

1040

1041

1042

1043



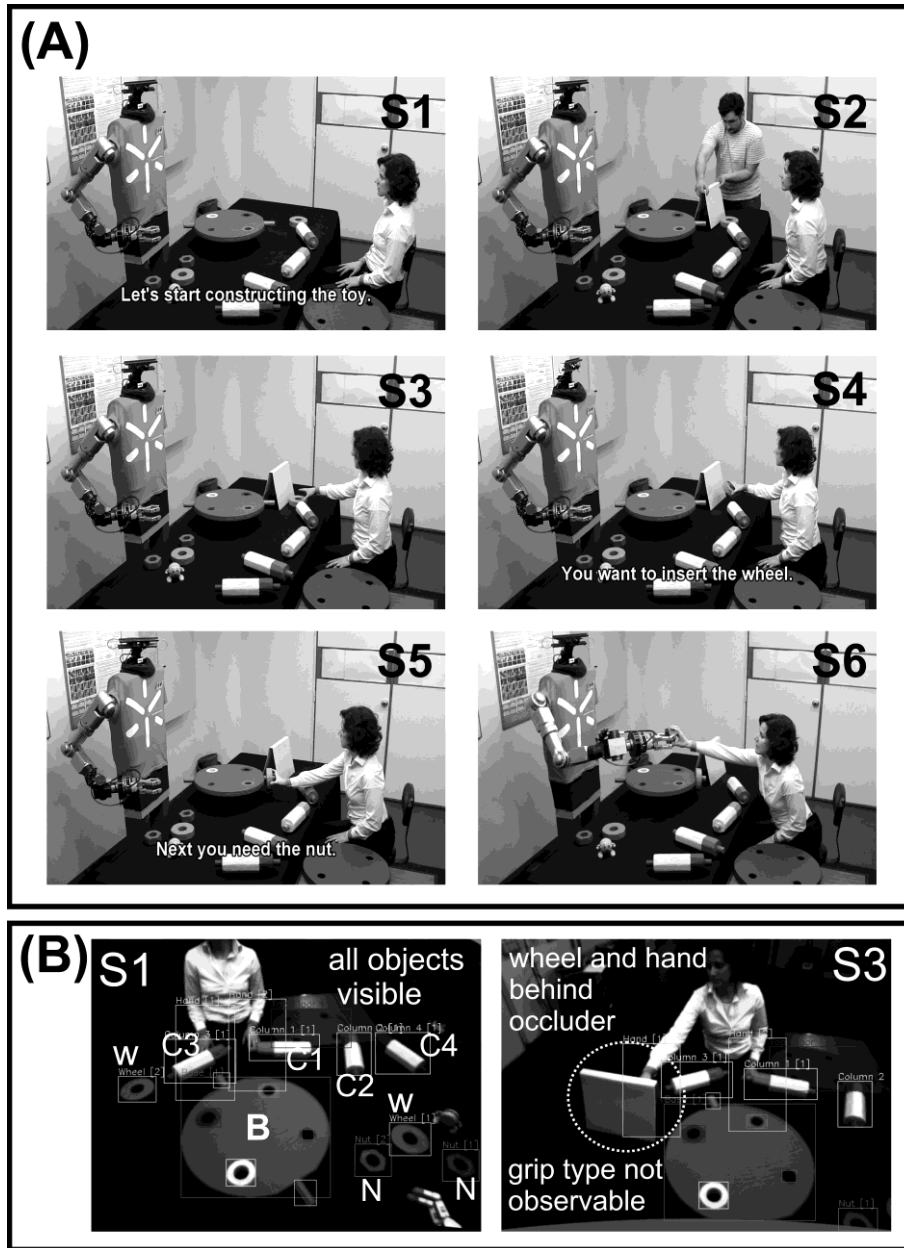
1044

1045

1046 **Figure 9:** Temporal evolution of field activity in the intention layer (IL)

1047 during time interval T1-T5 (video 1).

1048



1050

1051

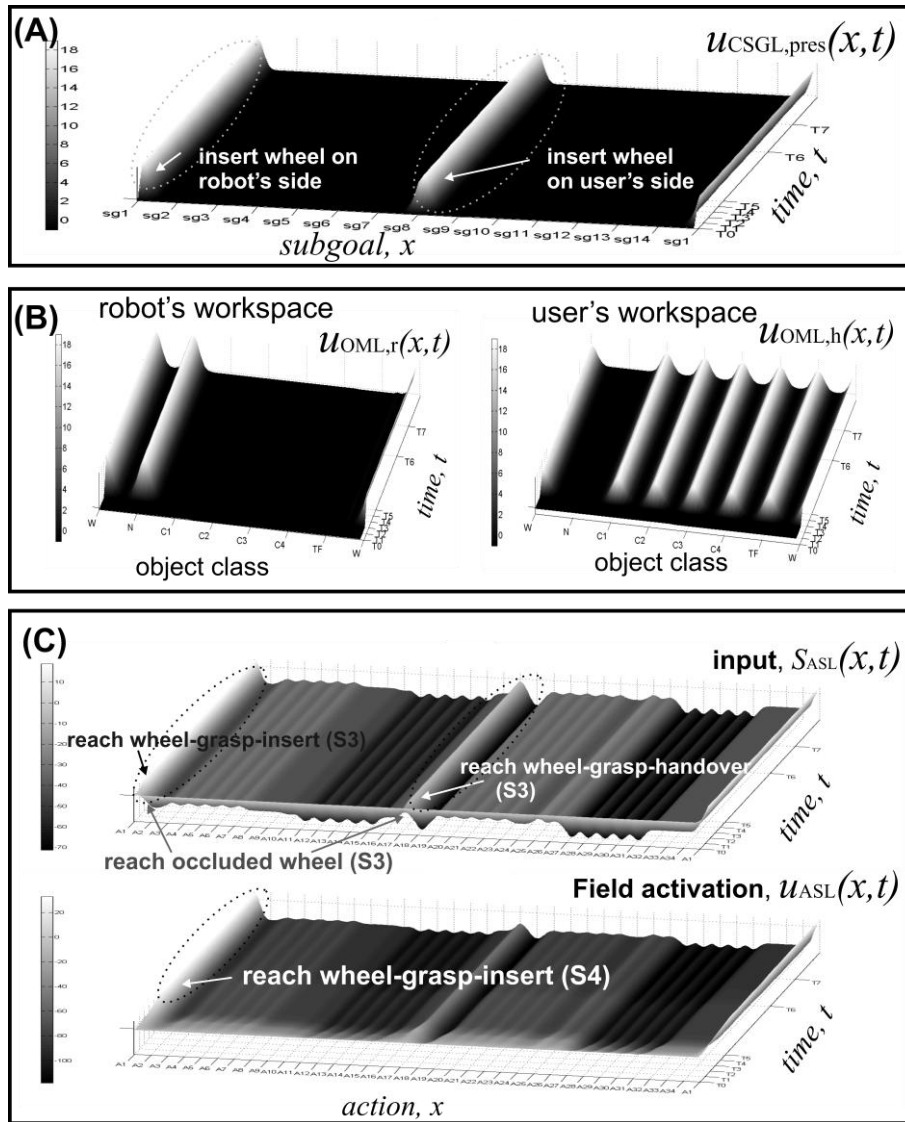
1052 **Figure 10:** Video 2: (A) Video snapshots showing action understanding of

1053 partially occluded actions (S1-S4) and anticipation of the user's future needs

1054 (S4-S6). (B) Snapshots of the robot's vision system.

1055

1056



1057

1058

1059 **Figure 11:** Video 2: (A) Temporal evolution of field activity representing

1060 present possible subgoals (in the CSGL). (B) Temporal evolutions of field

1061 activity in the OML. (C) Temporal evolutions of input to the ASL (top) and

1062 the field activity of the ASL (bottom).

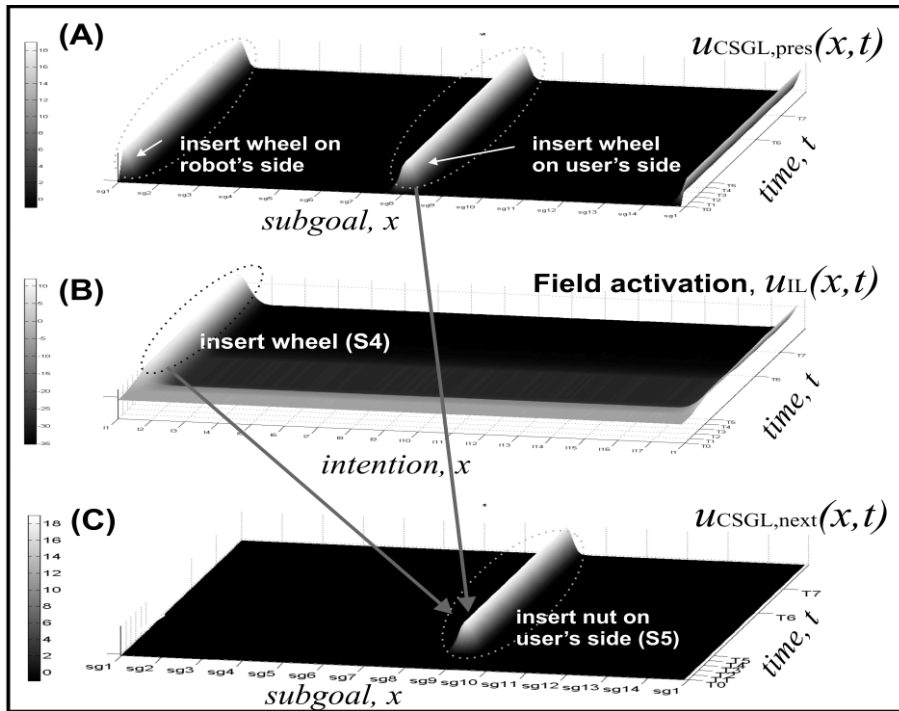
1063

1064

1065

1066

1067



1068

1069

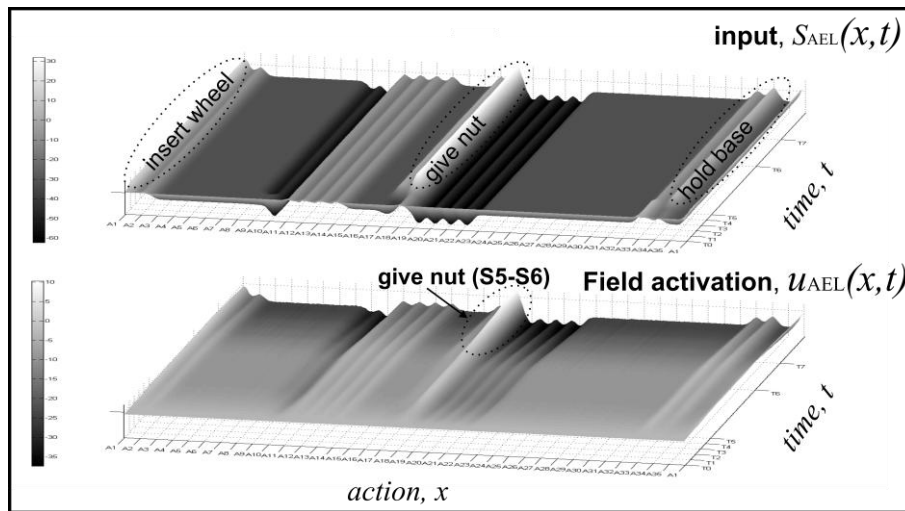
1070 **Figure 12:** Updating of the dynamic field representing subsequent subgoals
 1071 for the user based on a prediction of his/her current motor intention. Input
 1072 from the dynamic field encoding the current subgoals (A) and input from IL
 1073 (B) induces suprathreshold peak(s) of activation in the field encoding
 1074 subsequent assembly steps (C).

1075

1076

1077

1078



1079

1080

1081 **Figure 13:** Video 2: Anticipatory action selection. The temporal evolution
1082 of total input to AEL (top) and field activity in AEL (bottom) are shown.
1083 The robot decides to transfer a nut to the user (`give nut` action, snapshots
1084 S5-S6, Fig.10).

1085



1086

1087 **Figure 14:** Snapshots of Video 1 (long interaction scenario) in the time
 1088 interval T7-T11 are shown. They illustrate the robot's capacity to detect and
 1089 correct the user's intention errors and interpret a request gesture in a
 1090 context-sensitive manner (see the text for details).

1091

1092

1093

1094

1095

1096

1097

1098

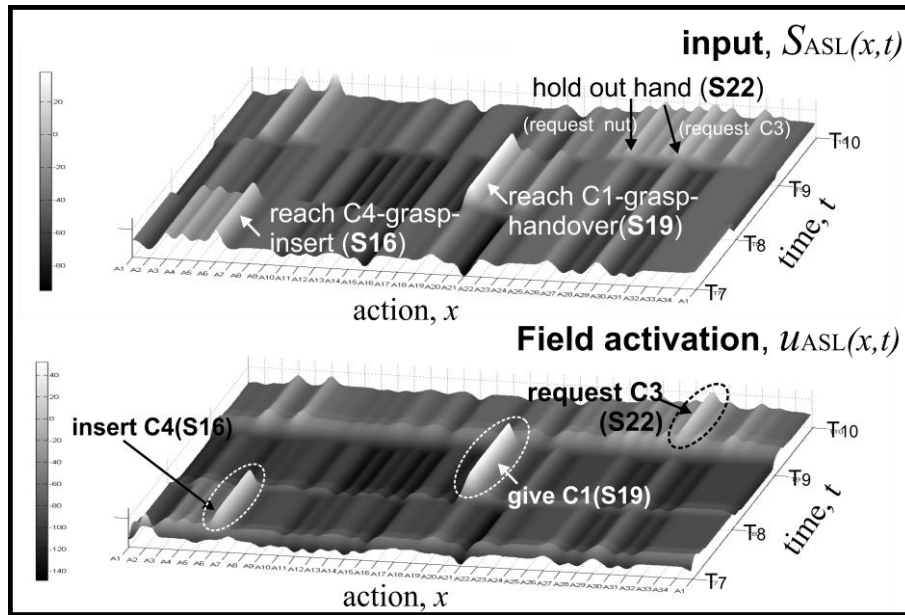
1099

1100

1101

1102

1103



1104

1105

1106 **Figure 15:** Error detection and context-sensitive interpretation of a request
 1107 gesture (video 1, time interval T1-T10). The temporal evolution of input to
 1108 the ASL (top) and the field activity in the ASL (bottom) are shown.

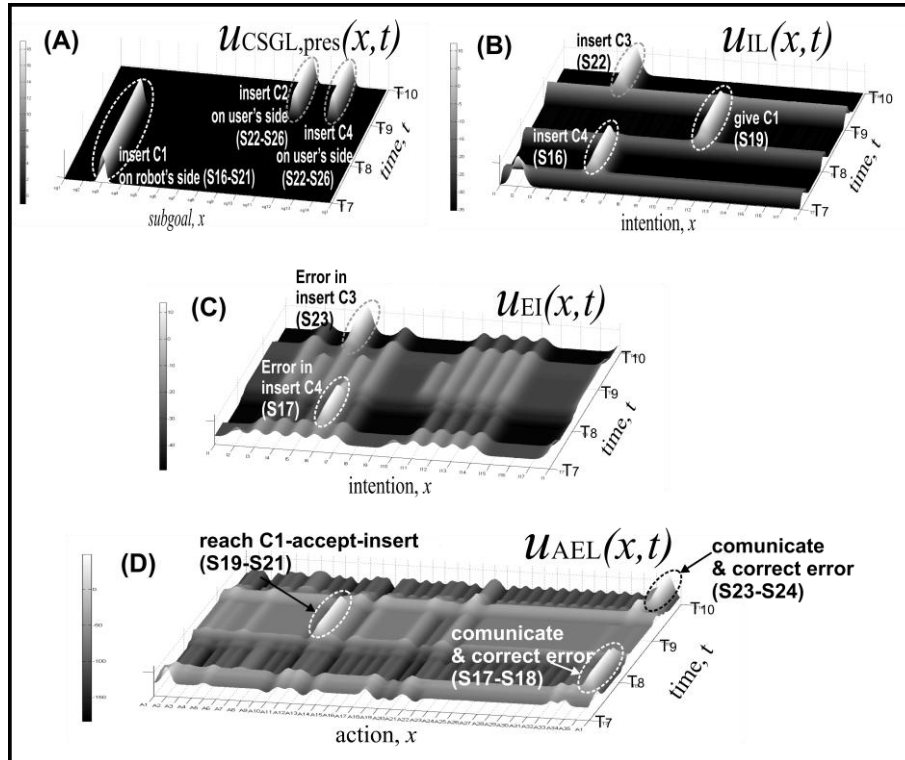
1109

1110

1111

1112

1113



1114

1115

1116 **Figure 16:** The temporal evolution of activity in different fields is shown
 1117 for the Video 1 (long interaction scenario) in the time interval T7-T10. (A)
 1118 Field of the CSGL representing currently available subgoals, (B) IL, (C) EI
 1119 encoding the errors in intention, (D) AEL.

1120

1121