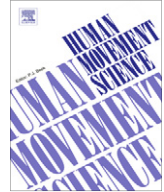




ELSEVIER

Contents lists available at ScienceDirect

Human Movement Science

journal homepage: www.elsevier.com/locate/humov

Neuro-cognitive mechanisms of decision making in joint action: A human–robot interaction study

Estela Bicho^{a,*}, Wolfram Erlhagen^b, Luis Louro^a, Eliana Costa e Silva^a

^aDepartment of Industrial Electronics, University of Minho, 4800-058 Guimarães, Portugal

^bDepartment of Mathematics and Applications, University of Minho, 4800-058 Guimarães, Portugal

ARTICLE INFO

Article history:

Available online xxx

PsycINFO Classification:

2340
3040
4140
4160

Keywords:

Motor planning
Decision making
Dynamic neural fields
Joint action

ABSTRACT

In this paper we present a model for action preparation and decision making in cooperative tasks that is inspired by recent experimental findings about the neuro-cognitive mechanisms supporting joint action in humans. It implements the coordination of actions and goals among the partners as a dynamic process that integrates contextual cues, shared task knowledge and predicted outcome of others' motor behavior. The control architecture is formalized by a system of coupled dynamic neural fields representing a distributed network of local but connected neural populations. Different pools of neurons encode task-relevant information about action means, task goals and context in the form of self-sustained activation patterns. These patterns are triggered by input from connected populations and evolve continuously in time under the influence of recurrent interactions. The dynamic model of joint action is evaluated in a task in which a robot and a human jointly construct a toy object. We show that the highly context sensitive mapping from action observation onto appropriate complementary actions allows coping with dynamically changing joint action situations.

© 2010 Elsevier B.V. All rights reserved.

1. Introduction

As an exquisitely social species humans are experts in cooperating with others when trying to achieve the goals of a common task (Sebanz, Bekkering, & Knoblich, 2006). In our everyday social interactions we continuously monitor the actions of our partners, interpret them in terms of their

* Corresponding author. Tel.: +351 253510181; fax: +351 253510189.

E-mail addresses: estela.bicho@dei.uminho.pt (E. Bicho), wolfram.erlhagen@mct.uminho.pt (W. Erlhagen), llouro@dei.uminho.pt (L. Louro), esilva@dei.uminho.pt (E. Costa e Silva).

outcomes and adapt our own motor behavior accordingly. Imagine for instance the joint action task of preparing a dinner table. The way a co-actor grasps a certain object (e.g., a coffee cup) or the context in which the motor act is executed (e.g., the cup may be empty or full) transmits to the observer important information about the co-actor's intention. Depending on the grip type, for instance, she/he may want to place the cup on the table or, alternatively, may have the intention to hand it over. Knowing what the other is going to do should facilitate motor programmes in the observer that serve the achievement of shared goals. Fluent and efficient coordination of actions among co-actors in a familiar task requires that the preparation of an adequate complementary action is a rather automatic and unconscious process. Since in sufficiently complex situations several possible complementary behaviors may exist this process necessarily includes a decision-making operation.

The long-term goal of our research group is to build robots that are able to interact with users in the same way as humans interact with each other in common tasks. Our research strategy to achieve this challenging objective is to develop and test control architectures that are strongly inspired by neuro-cognitive mechanisms underlying human joint action. We believe that implementing a human-like interaction model in an autonomous robot will greatly increase the user's acceptance to work with an artificial agent since the co-actors will become more predictable for each other (for a survey of challenges for socially interactive robots see [Fong, Nourbakhsh, and Dautenhahn \(2003\)](#)). Such an interdisciplinary approach constitutes not only a promising line of research towards human-centered robots but also offers unique possibilities for researchers from neuroscience and cognitive science. Synthesizing cooperative behavior in an artificial but naturally inspired cognitive system allows them in principle to test their theories and hypothesis about the mechanisms supporting social interactions ([Dominey & Warneken, in press](#)).

The focus of this paper is on flexible action planning and decision formation in cooperative human-robot interactions that take into account the inferred goal of the co-actor and other task constraints. An impressive range of experimental evidence accumulated over the last two decades supports the notion that a close perception-action linkage provides a basic mechanism for real-time social interactions ([Newman-Norlund, Noordzij, Meulenbroek, & Bekkering, 2007](#); [Wilson & Knoblich, 2005](#)). A key idea is that action observation leads to an automatic activation of motor representations that are associated with the execution of the observed action. It has been suggested that this resonance of motor structures supports an action understanding capacity ([Blakemore & Decety, 2001](#); [Fogassi et al., 2005](#); [Rizzolatti, Fogassi, & Gallese, 2001](#)). By internally simulating action consequences using his own motor repertoire the observer may predict the consequences of others' actions. Direct physiological evidence for such a perception-action matching system came with the discovery of the mirror neurons first described in premotor cortex of macaque monkey (for a review see [Rizzolatti and Craighero \(2004\)](#)). Mirror neurons are a particular class of visuomotor neurons that are active both during the observation of goal-directed actions such as reaching, grasping, holding or placing an object and during the execution of the same class of actions. Although action understanding is the dominant hypothesis about the functional role of the motor resonance mechanism it has been suggested that it may also contribute to motor planning and action preparation. Typically it is assumed that a direct activation of the corresponding motor program explains the evidence found in many behavioral experiments for a tendency of an automatic imitation of observed actions (e.g., [Brass, Bekkering, & Prinz, 2001](#), for a review see [Wilson & Knoblich, 2005](#)). Such a tendency is of course not beneficial for cooperative joint action which normally requires the facilitation of a complementary motor behavior. Recent findings in neuroimaging and behavioral studies provide evidence however that goal and context representations may link an observed action to a different but functionally related motor response ([Newman-Norlund, van Schie, van Zuijlen, & Bekkering, 2007](#); [van Schie, van Waterschoot, & Bekkering, 2008](#)). These studies clearly demonstrate that the mapping between action observation and action execution is much more flexible than previously thought.

Here we present a dynamic model that implements such a flexible perception-action linkage as a means to achieve an efficient coordination of actions and decisions between co-actors in a joint action task. We report results of our ongoing evaluation of the model as part of the control architecture of a humanoid robot that assembles together with a human user a toy object from its components ([Bicho, Louro, Hipolito, & Erhagen, 2008, 2009](#)).

The model is based on the theoretical framework of Dynamic Neural Fields (DNFs) that has been originally proposed to explain the firing patterns of neural populations in the cortex (Amari, 1977; Wilson & Cowan, 1973; see also Grossberg (1973) for a related approach). The architecture of this model family reflects the hypothesis that strong recurrent interactions in local pools of neurons form a basic mechanism of cortical information processing. DNFs have been first introduced into the motor domain as neuro-inspired models of sensorimotor decisions in simple reaching and saccadic eye movement tasks (Erlhagen & Schöner, 2002; Schöner, Kopecz, & Erlhagen, 1997; Wilmzig, Schneider, & Schöner, 2006). In these applications, the dynamic fields represent parameters of the movement such as for instance extent and direction. The neural activation patterns encoding these parameters evolve continuously in time under the influence of inputs representing sensory evidence and prior task knowledge but are mainly shaped by the interplay of local-excitatory and long-range inhibitory interactions within the population. Due to the recurrent interactions, the patterns may become self-stabilized in the absence of any external input. Such stable states of the field dynamics reflect decisions between multiple movement alternatives since a competition process mediated by lateral inhibition leads to a suppression of activation in neural pools that are less supported by input from external sources.

Here we apply the basic concepts of DNF-models to address the problem of selecting a specific action among a set of possible complementary behaviors. The competing neural populations in the model thus encode entire object-directed motor acts like reaching, grasping, placing or combinations of these motor primitives. The high level of abstraction of the neural representations fits well to the fundamental properties of most mirror neurons that encode the goal of an action (e.g., the effector interacts with the object in an efficient way) independent of the fine details of the movement kinematics (Rizzolatti & Craighero, 2004). Of particular interest is a class of mirror neurons that reveals a broad matching between action observation and execution (e.g., involving different effector and/or postures), which could in principle support a flexible perception-action coupling in cooperative settings (Newman-Norlund, van Schie et al., 2007).

The DNF-model of joint action extends our previous modeling work on action understanding and goal-directed imitation (Erlhagen, Mukovskiy, & Bicho, 2006a, 2006b; Erlhagen, Mukovskiy, Chersi, & Bicho, 2007). It consists of a multi-layered network of reciprocally connected neural populations that represent in their activation patterns specific task-relevant information. Decision making in the joint assembly task is implemented as a dynamic process that continuously integrates over time information about the inferred goal of the co-actor (obtained through motor simulation), shared knowledge about what the two actors should do (construction plan), and contextual information (e.g., the spatial distribution of objects in the working area). For generating the overt motor behavior of the robot in the joint construction scenario, we apply a posture-based motor planning and execution model. It translates the abstract decision about a complementary action (e.g., grasping an object to hold it out for the partner) into a realistic, collision-free trajectory.

The paper is organized as follows: Section 2 introduces the joint construction task and briefly describes the robotic platform used for the human-robot interaction experiments. Section 3 gives an overview of the motor planning model that we have used to generate the overt behavior of the robotic arm and hand. Section 4 gives an overview of the model for decision making in joint action. Section 5 presents the basic concepts of the Dynamic Neural Field framework and summarizes the DNF-based implementation of the joint action model. The selection of appropriate complementary actions in different joint action contexts are described in the results Section 6. The paper ends with a discussion of results, concepts and future research.

2. Joint construction task

To validate the dynamic field model of joint action we have chosen a task in which a robot collaborates with a human in constructing a toy object from components that are initially distributed on a table (Fig. 1).

The task requires only a limited number of different motor behaviors to be performed by the team but is complex enough to show the impact of goal inference, shared task knowledge and context on

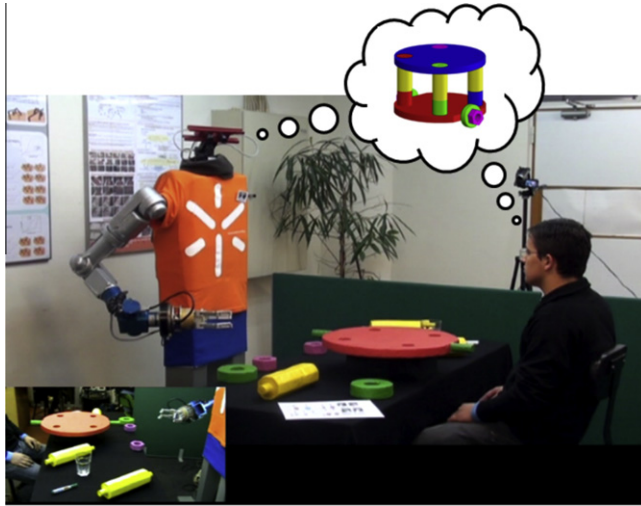


Fig. 1. Human–robot joint action scenario. The task for the team is the joint construction of a toy object.

action selection. The components that have to be manipulated by the robot were designed to limit the workload for the vision and the motor system of the robot. The toy object consists of a round platform with an axle on which two wheels have to be attached and fixed with a nut. Subsequently, 4 columns have to be plugged into holes in the platform. The placing of another round object on top of the columns finishes the task. It is assumed that each teammate is responsible to assemble one side of the toy. Since the working areas of the human and the robot do not overlap, the spatial distribution of components on the table obliges the team to coordinate in addition to handing-over sequences. It is further assumed that both partners know the construction plan and keep track of the subtasks which have been already completed by the team. As part of the dynamic control architecture, the plan is implemented based on the concepts of the dynamic field theory. Hand-designed connections between populations encoding subsequent subtasks define the logical order of the assembly work (for more details see Section 4). Since the desired end state does not uniquely define this order, at each stage of the construction the execution of several subtasks may be simultaneously possible. The main challenge for the team is thus to efficiently coordinate in space and time the decision about actions to be performed by each of the co-actors.

For the experiments we used a robot built in our lab (Silva, Bicho, & Erlhagen, 2008). It consists of a stationary torus on which an arm with 7 degrees of freedom and a 3-fingered hand, and a stereo camera head are mounted. A speech synthesizer allows the robot to communicate the result of its reasoning and decision processes to the human user. The information about object type, position and pose as well as about the state of the construction is provided by the camera system. The object recognition combines color-based segmentation with template matching derived from earlier learning examples (Westphal, von der Malsburg, & Würtz, 2008). The same technique is also used for the classification of object-directed, static hand postures such as grasping and communicative gestures such as pointing or demanding an object.

3. Movement planning

For the human–robot experiments, a decision of the robot to execute a specific complementary action has to be translated into a fluent, smooth and collision-free arm trajectory. A complementary behavior may consist of a simple pointing towards a target object, but may also involve more complex movements such as grasping a component with a specific grip or attaching different components of the toy object to each other. These goal-directed movements define the motor repertoire of the robot

in the construction task. To generate complete temporal motor behaviors of the robotic arm and hand we use an approach that is inspired by the posture model of Rosenbaum and colleagues (Meulenbroek, Rosenbaum, Jansen, Vaughan, & Vogtn, 2001; Rosenbaum, Meulenbroek, Vaughan, & Jansen, 2001). This model has been proven to generate different types of realistic movements such as reaching, grasping and manipulation of objects, and presents an elegant obstacle avoidance mechanism. The posture model was first introduced for planar movements and was recently extended to the planning in a 3D workspace (Vaughan, Rosenbaum, & Meulenbroek, 2006). A key assumption is that the planning of movements in joint space can be divided into two sub-problems: end posture selection and trajectory selection. Here we give an overview of the implementation of this two-step planning process in the robot. The model has been described with more technical details elsewhere (Costa e Silva, Bicho, Erlhagen, & Meulenbroek, submitted). The planning system first selects a goal posture from the set of all postures that (1) allows an object to be grasped without collision with any obstacle, and (2) minimizes the displacement costs from the beginning to the end of the movement. Mathematically, the selection process can be formalized as a nonlinear constraint optimization problem. It is numerically solved taking into account the information about object type, position and orientation (provided by the vision system), as well as the information represented by the activation patterns in the dynamic field model about grip type and hand orientation relative to the object. Subsequently, the trajectory is selected by computing for each of the 10 joints of the robotics arm and hand its trajectory, i.e., time history of position, velocity and acceleration, from initial to end posture. Since the minimum jerk principle is applied (Flash & Hogan, 1985), the movements of the joint follow a bell-shaped velocity profile, resulting in a smooth straight-line movement in joint space. This joint trajectory defines the direct movement without checking if an obstacle blocks a certain area in posture space. To detect potential collision with an intermediate obstacle, the object to be grasped or a target object, the planning system uses direct (forward) kinematics to internally simulate movement execution from start to end. If no collision is anticipated, the movement is executed, otherwise the system searches for a feasible movement. For finding this alternative a suitable bounce posture is selected. This bounce posture serves as a subgoal for a back-and-forth movement, which is superimposed on the direct movement. The end posture that is finally reached is the same as for the direct movement, only the selected path differs to guarantee collision avoidance. The bounce posture is found by solving a similar constrained optimization problem as applied for the end posture. It minimizes the displacement of the joints. To generate the movement, the desired joint position and time interval, given by the planning model, is sent to the low-level arm and hand controllers, using the high-level interface functions provided by the manufacturer (AMTEC/SCHUNK and BARRETT Technology, respectively). They guarantee that the planned trajectory is executed in the desired time interval. The real-time interaction experiments with human users show that the movements of the robot are perceived as smooth and goal-directed but slower compared to human motion. Moreover, a direct comparison with human data in reach to grasp tasks reveals that the generated arm and hand trajectories reflect several characteristics observed in biological motion such as for instance a biphasic tangential velocity profile or a maximum grip aperture that occurs during the second half of the movement (Lommertzen, Costa e Silva, Cuijpers, & Meulenbroek, 2008).

4. Model overview

Fig. 2 presents a sketch of the DNF-based architecture for decision making in cooperative joint action. It consists of various layers each containing one or more neural populations encoding information specific to the construction task (a detailed description of the labels in each layer is given in the [Supplemental material](#)). The lines indicate the connectivity between individual populations in the network. Basically, the architecture implements a flexible mapping from an observed action of the co-actor onto a complementary motor behavior.

The multi-layered architecture extends a previous DNF-model of the STS-PF-F5 mirror circuit of monkey (Erlhagen et al., 2006a) that is believed to represent the neural basis for a matching between the visual description of an action in area STS and its motor representation in area F5 (Rizzolatti & Craighero, 2004). This circuit supports a direct and automatic imitation of an action performed by

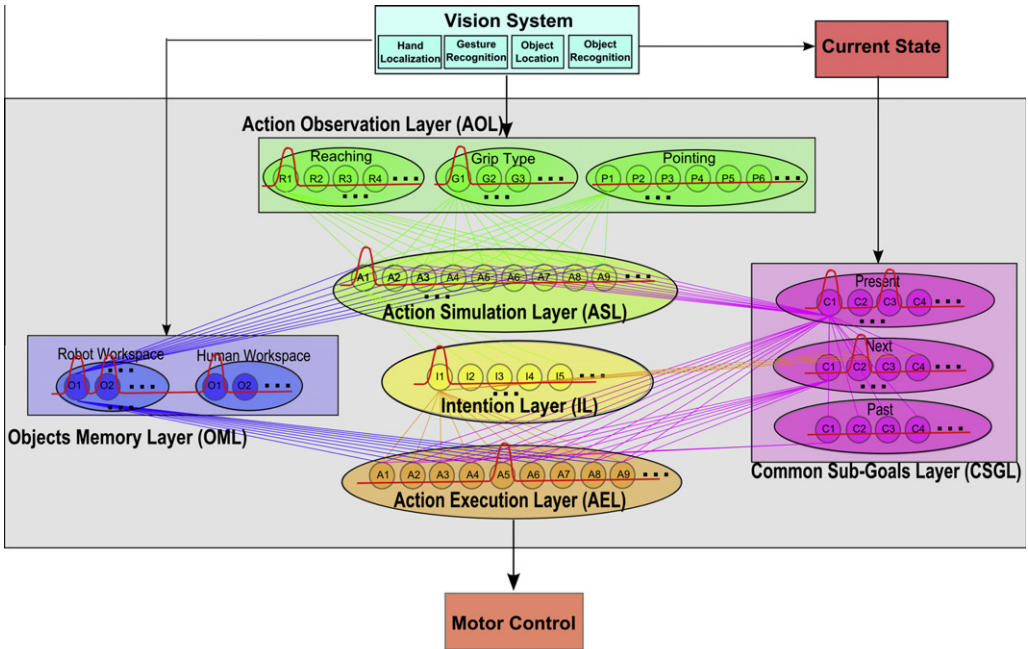


Fig. 2. Joint action model. It implements a flexible mapping from observed actions (layer AOL) onto complementary actions (layer AEL) taking into account the inferred action goal of the partner (layer IL), contextual cues (layer OML) and shared task knowledge (layer CSGL). The goal inference capacity is based on motor simulation (layer ASL).

another individual. Importantly for joint action, however, the model allows also for a flexible perception-action coupling by exploiting the existence of object-directed action chains in the middle layer PF (Fogassi et al., 2005) that are linked to the representations of their final goals or outcomes in prefrontal cortex (PFC). The automatic activation of a particular chain during action observation (e.g., reaching-grasping-placing) drives the connected goal representation which in turn may bias the decision processes in layer F5 towards the selection of a complementary rather than an imitative action. Consistent with this model prediction, a class of mirror neurons has been reported in F5 for which the effective observed and effective executed actions are logically related (e.g., implementing a matching between placing an object on the table and bringing the object to the mouth (di Pellegrino, Fadiga, Fogassi, Gallesse, & Rizzolatti, 1992)). For the robotics work we refer to the three layers of the matching system as the action observation (AOL), action simulation (ASL) and action execution (AEL) layers, respectively. An observed object-related hand movement that is recognized by the vision system as a particular primitive is represented in AOL. In the action simulation layer (ASL) populations encode entire chains of action primitives that are in the motor repertoire of the robot (e.g., reaching-grasping-placing/attaching a particular part) or communicative hand gestures (e.g., pointing towards or requesting a part). They are connected to population representations of the associated end states or goals in the intention layer IL (e.g., right wheel attached). The activation of a particular chain during action observation thus allows the robot to predict the co-actor's motor intentions by internally simulating the action outcomes. Very often, however, the observation of a particular motor act alone (e.g., grasping) is not sufficient to make this prediction since the motor act may be part of several chains. To solve this ambiguity, the neural populations in ASL get additional inputs from connected populations representing situational context and/or prior task knowledge about what the co-actor should do in a particular situation. An important contextual cue is the spatial distribution in the workspace of parts necessary for the assembly work. The object memory layer (OML) encodes memorized information about the

position of these parts in each of the two working areas, separately for each object type. The common subgoals layer (CSGL) encodes the currently available subgoals as well as the subtasks that have been already accomplished by the team. The available subgoals are continuously updated based on feedback from the vision system in accordance with the construction plan. The information about the sequential order in which subtasks have to be accomplished (e.g., attach right wheel first and subsequently fix it with a nut) is encoded in the synaptic links between populations representing these subgoals in two different neural fields (indicated with labels 'present' and 'next' in layer CSGL, see Fig. 2). Input from the vision system signaling the achievement of a certain subtask activates the respective population representation in the first layer which in turn drives through the connections the populations representing the next possible assembly steps in the second layer. To guarantee pro-active behavior of the robot, the model implements the possibility to update the current subgoals also based on input from IL representing the predicted motor intention of the co-actor. This allows the robot to start preparing an action serving a subsequent goal (e.g., transferring a nut to the co-actor for fixing the wheel) ahead of the realization of the preceding subtask (e.g., co-actor is going to attach right wheel). In the action execution layer (AEL) populations that encode the same action sequences and communicative gestures like the ASL compete for expression in overt behavior. They integrate input from the IL, OML and CSGL.

To give an example of the dynamic decision making process implemented in the field architecture, think of the situation that the co-actor reaches towards a wheel in his working area. The wheel on his construction side has been already attached, but not the wheel on the side of the observer. The available information about active and already accomplished subtasks together with the observed hand motion activates automatically the chain representation of a 'reach-grasp wheel-handover' behavior in the ASL which subsequently activates the motor intention 'handover wheel' in the IL. As a consequence, the robot may prepare at the time of the reaching for receiving the wheel. Now imagine that the same motor act is observed at the start of the construction, that is, the wheel on the co-actor's side has not been attached yet. Consequently, specific input from CSGL in support of the object transfer hypothesis is missing. Now input from the AOL representing the type of the observed grasping behavior (top versus side grip) may decide which of the two possible chains associated with different motor intentions may become activated. A top grip (i.e., a grip from above) is usually used for directly attaching the wheel, whereas grasping from the side is the most secure way to hand over the wheel to the partner. In the latter case, the robot will prepare a complementary grasping behavior to receive the wheel. In the former case, an adequate complementary behavior of the robot might be to reach for a wheel in its workspace to attach it on its construction side.

5. Model details: Basic concepts of dynamic neural field theory

Each layer of the model is formalized by one or more Dynamic Neural Fields (DNFs). The DNFs implement the idea that task-relevant information about action goals, action primitives or context is encoded by means of activation patterns of local populations of neurons. In the action observation layer for instance, high levels of activation of a neural pool representing a certain grasping behavior means that the specific grip type has been detected and classified by the vision system whereas a low activation level indicates that information about the specific grip type is currently not processed (see Fig. 3, top).

As shown in the bottom panel of Fig. 3, the activation patterns are initially triggered by input from connected populations and/or sources external to the network like for instance the vision system. In the example, the input to the population representing an above grip is stronger compared to the input to the population representing a side grip.

The activation patterns may become self-sustained in the absence of any input due to the recurrent interactions within the local populations. To guarantee the existence of self-stabilized solutions of the field dynamics, the pattern of recurrent interactions between cells must be spatially structured with excitation dominating at small and inhibition at larger distances. The distance between neurons may be defined in anatomical space (as in the original work by Wilson and Cowan (1973) and Amari (1977)) or in a more abstract space given by some feature dimension that the neurons encode (e.g.,

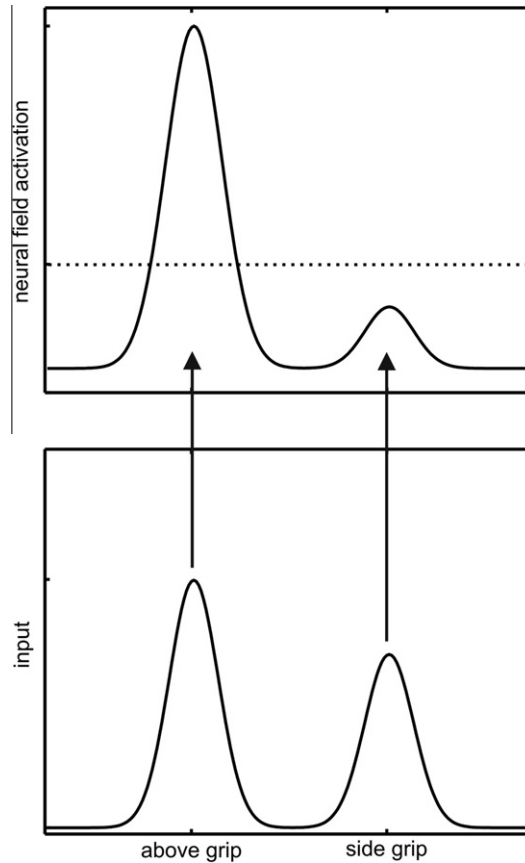


Fig. 3. Top panel: Representation of information in the dynamic field framework for the example of a grasping behavior (top panel). Activation patterns may represent through a localized peak of activation with an excitation level above a certain threshold (dotted line) the presence of information about a certain grasping behavior (above grip). Low excitation levels indicate that information about the alternative grasping behavior (side grip) is currently not processed. Bottom panel: The activation patterns of the neural populations are initially triggered by external input which may come from connected populations or like in the example from the vision system.

movement direction or direction in visual space). For high-dimensional spaces representing for instance different grasping behaviors the metric distance is not directly observable. However, it may still be defined operationally by the degree of overlap between their neural representations. For functionally distinct behaviors associated with the achievement of different goals, we assume that they are represented by separate neural sub-populations (all able to self-sustain high levels of activation) located at a distance that guarantees a purely inhibitory interaction between the neurons of these pools.

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

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

where $\tau_i > 0$ and $h_i > 0$ define the time scale and the resting level of the field dynamics, respectively. The integral term describes the intra-field interactions. It is assumed that the interaction strength, $w_i(x, x')$, between any two neurons x and x' depends only on the distance between field locations.

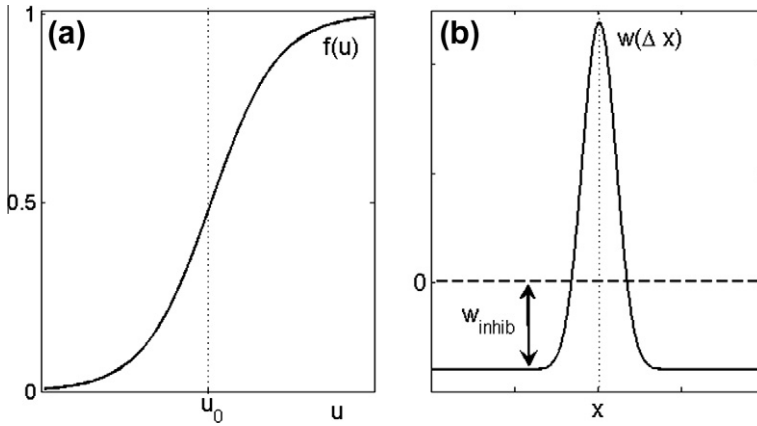


Fig. 4. (a) Nonlinear output function $f(u)$ of sigmoid shape with threshold u_0 . (b) Synaptic weight function $w(\Delta x)$, $\Delta x = x - x'$, of center-surround type. The recurrent interactions are excitatory for neighboring neurons, whereas for larger distances inhibition dominates.

For the present implementation we used the following weight function of lateral-inhibition type (Fig. 4b):

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

where $A_i > 0$ and $\sigma_i > 0$ describe the amplitude and the standard deviation of a Gaussian, respectively. For simplicity, the inhibition is assumed to be constant, $w_{inhib,i} > 0$, meaning that integration of inhibitory input does not change with distance between field sites. Only sufficiently activated neurons contribute to interaction. The threshold function $f_i(u_i)$ is chosen of sigmoidal shape with slope parameter β and threshold u_0 (Fig. 4a):

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

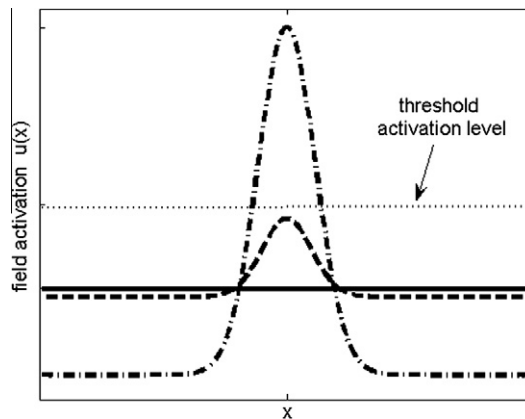


Fig. 5. Bi-stable system: In response to weak input the dynamic field generates an input driven activation pattern (dashed line). When the input drives the field activation beyond a threshold value the homogeneous state (solid line) loses stability and a self-stabilized activation pattern evolves (dashed-dotted line).

As illustrated in Fig. 5, the model parameters are adjusted to guarantee that the field dynamics is bi-stable (Amari, 1977), that is, the attractor state of a self-stabilized activation pattern (dashed-dotted line) coexists with a stable homogenous activation distribution (solid line) that represents the absence of specific information (resting level). If the summed input, $S_i(x, t)$, to a local population is sufficiently strong to drive the field activation beyond a certain threshold, the homogeneous state loses stability and a localized pattern in the dynamic field evolves. Weaker external signals lead to a subthreshold, input-driven activation pattern (dashed line) in which the contribution of the interactions is negligible.

Normally, a constant input from a single population does not drive directly connected populations. It may play nevertheless an important role for the processing in the joint action circuit. The preshaping by weak input brings populations closer to the threshold for triggering the self-sustaining interactions and thus biases the decision processes linked to behavior. Much like prior distributions in the Bayesian sense, multi-modal patterns of subthreshold activation in for instance the action execution layer (AEL) may represent the probability of different complementary actions (Erlhagen & Schöner, 2002).

The summed input from connected fields u_i is given as $S_i(x, t) = k \sum_j S_j(x, t)$. The parameter k scales the total input relative to the threshold for triggering a self-sustained pattern. This guarantees that the inter-field coupling is weak and the field dynamics is dominated by the recurrent interactions. The input from each connected field u_i is modeled by Gaussian functions. As shown in Fig. 6, the input from a connected population j in layer u_i to a target population m in layer u_i is modeled by a Gaussian function. This input is applied whenever the activation in population j is beyond the threshold for a self-stabilized activation peak. The total input from all sub-populations in field u_i to field u_i is mathematically described by:

$$S_i(x, t) = \sum_m \sum_j a_{mj} c_{ij}(t) \exp(-(x - x_m)^2 / 2\sigma^2). \quad (4)$$

Here $c_{ij}(t)$ is a function that signals the existence or evolution of a self-sustained activation pattern in field u_i centered at position y_j (i.e., signals that subpopulation j in u_i is active), and a_{mj} is the inter-field synaptic connection between subpopulation j in u_i to subpopulation m in u_i .

The existence of a self-stabilized activation peak in a dynamic field is closely linked to decision making. In layers ASL, IL and AEL subpopulations encoding different chains (ASL), goals (IL) and complementary actions (AEL), respectively, interact through lateral inhibition. Fig. 7 shows the temporal evolution of two competing populations encoding two different actions. The inhibitory interactions cause the suppression of activity below resting level in competing neural pools whenever a certain subpopulation becomes activated above threshold in response to external input.

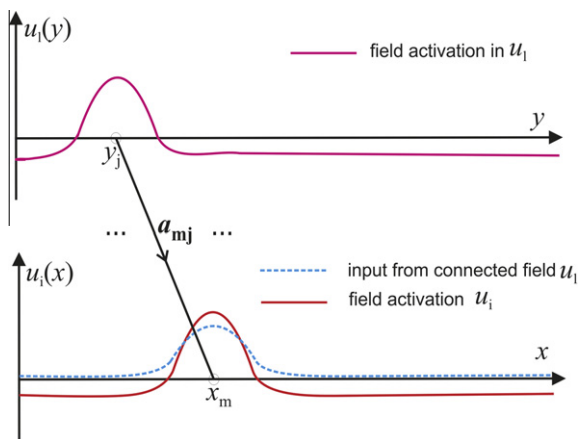


Fig. 6. Schematic view of the input from a population j in layer u_i that appears to be activated beyond threshold level to a target population m in u_i . The input is modeled as a Gaussian function.

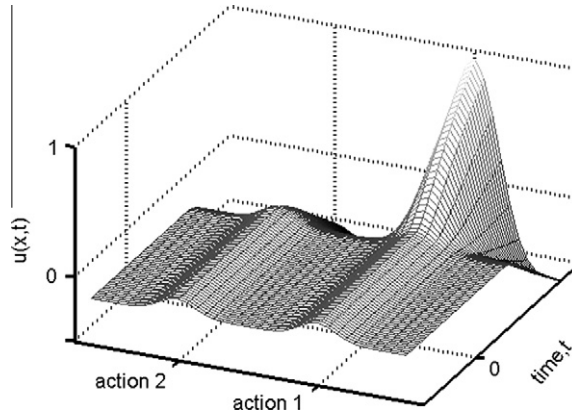


Fig. 7. Decision making in neural fields. The temporal evolution of activity in two populations representing two different actions is shown. An input of equal strength is applied at time $t = 0$ to both populations. Note that at that time the neurons appear to be already pre-activated above resting level due to subthreshold input from other populations. The population representing action 1 wins the competition since its level of pre-activation is slightly higher than the level of the neural pool representing action 2.

The attractor state of a self-sustained activation peak can be used to implement a working memory function, which happens for sufficiently small values of the global inhibitory parameter $h_i > 0$. Conversely, for sufficiently large values of $h_i > 0$ the existing suprathreshold self-sustained activation becomes unstable and field activation decays back to resting level (the field becomes mono-stable). To implement these working memory versus forgetting mechanisms we defined a proper dynamics for the global inhibitory parameter $h_i > 0$ (Bicho, Mallet, & Schöner, 2000):

$$\frac{dh_i}{dt} = -\left(\frac{1}{\tau_{i,\min}}\right)(1 - c_h)(h_i - h_{i,\min}) - \left(\frac{1}{\tau_{i,\max}}\right)c_h(h_i - h_{i,\max}) \quad (5)$$

This dynamics increases h_i toward $h_{i,\max}$ (destabilizing memory), at a time scale $\tau_{i,\max}$, while a suprathreshold pattern of activation exists ($c_h = 1$). It restores h_i to $h_{i,\min}$ (to enable the working memory), at a time scale $\tau_{i,\min}$. The restoring process is much faster than the destabilizing process ($\tau_{i,\min} \ll \tau_{i,\max}$), so that quickly after forgetting, the field dynamics is again able to self-sustain (i.e. memorize) a suprathreshold pattern of activation. The presence of a self-sustained peak is signaled by

$$c_h = [H(N_u) - H(N_s)]H(N_u) \quad (6)$$

where $H(\cdot)$ is a Heaviside step function, and

$$N_u = \int H(u(x))dx \quad (7)$$

is the total suprathreshold activation in the field and

$$N_s = \int H(S(x))dx \quad (8)$$

is the total external input activation to the field.

In the OML and the CSGL self-sustained peaks encode memorized information about the location of the relevant parts in the two working areas and the already achieved and currently available subgoals, respectively. Since multiple potential target objects and subgoals may exist simultaneously, the field dynamics must support multi-peak solutions. Although specific lateral interaction functions exist that support such solutions (Laing, Tray, Gutkin, & Ermentrout, 2002), we implemented for simplicity synaptic weight functions, in these layers, with limited spatial ranges of lateral inhibition to avoid competition (i=CSGL, OML):

$$w_i(x - x') = A_{i,exc} \exp(-(x - x')^2 / 2\sigma_{i,exc}^2) - A_{i,ini} \exp(-(x - x')^2 / 2\sigma_{i,ini}^2) \quad (9)$$

with $A_{i,exc} > A_{i,ini}$ and $\sigma_{i,exc} > \sigma_{i,ini}$.

6. Results

In the following we present results of real-time human–robot interactions in the joint construction scenario. It is assumed that both actors know the construction plan and are able to perceive the state of the construction. Each actor thus knows in principle what assembly steps the team still has to perform. The examples are chosen to illustrate the impact of action observation and context on action selection from the perspective of the robot. In each case, the initial spatial distribution of objects in the two working areas obliges the robot to continuously monitor and interpret the actions of its human partner, since actions such as handing-over or demanding objects are necessarily involved. A detailed discussion of the goal inference capacity is not within the scope of the present paper. We therefore focus on examples at the beginning of the joint assembly work for which action understanding is straightforward (for details of the motor simulation mechanism see Erlhagen et al. (2006a) and Erlhagen, Mukovskiy, Chersi, & Bicho (2007)).

The connections between the neural populations in the various fields are hand-coded meaning that the different inputs that may bias the selection of a particular complementary action are pre-defined. The robot shows flexibility in its behavior since in the complex dynamical system of interacting populations the decision for a certain action depends on both the informational content of the various inputs to populations in AEL and their timing. Changes in the time course of activity in a connected field due to competition between neural pools or noisy input data for instance may thus affect which complementary behavior the robot selects.

As summarized in Table 1, the total number of goal-directed sequences and communicative gestures that represent relevant complementary behaviors in AEL is restricted to 9 alternatives. At any point of time of the human–robot interaction only a few of these alternatives are simultaneously possible, that is, are supported by input from connected layers. It is important to stress, however, that the dynamic decision making implemented in AEL also works for more complex situations with a larger set of possible response alternatives (e.g., at later stages of the construction process). In line with the classical Hick–Hyman law, the number of alternatives only affects the time it takes to stabilize a peak solution representing a decision in the dynamic neural field (Erlhagen & Schöner, 2002).

Videos of the human–robot experiments can be found at <http://dei-s1.dei.uminho.pt/pessoas/este-laj/ASTvideos.htm>. The robot uses speech to communicate the outcome of the goal inference and decision making processes implemented in the dynamic field model to the human co-actor. As our studies with naive users show, this basic form of verbal communication facilitates natural and fluent interaction with the robot (Bicho, Louro, & Erlhagen, 2010).

Numerical values for the dynamic field parameters and synaptic inter-field connections may be found in the Supplemental material.

Table 1

Goal-directed sequences and communicative gestures.

Action	Sequence of motor primitives	Short description
A_1	Reach wheel → grasp → attach	Attach wheel
A_2	Reach wheel → grasp → handover	Give wheel
A_3	Reach hand → grasp wheel → attach	Receive wheel to attach
A_4	Reach nut → grasp → attach	Attach nut
A_5	Reach nut → grasp → handover	Give nut
A_6	Reach hand → grasp nut → attach	Receive nut to attach
A_7	Hold out hand	Request piece
A_8	Point to wheel	Point to wheel
A_9	Point to nut	Point to nut

6.1. Same observed action, different complementary behaviors

According to the plan, the construction starts with attaching the two wheels and subsequently fixing each of them with a nut. Fig. 8 (video snapshots), Fig. 9 (DNF in AEL) and Fig. 10 (DNF in IL), show that the same observed gesture may have a different interpretation depending in the context in which it occurs and thus leads to a different complementary action.

In this experiment the two wheels and one nut are initially located in the workspace of the robot. The robot may thus decide to choose an action that serves one of two possible subtasks: attaching the wheel on its side (A_1) or transferring a wheel to the partner (A_2) so that he can attach it on his side. As shown in snapshots S2–S3 the co-actor first reaches his empty hand towards the robot.

The robot interprets this gesture as a request for a wheel since attaching a wheel on the side of the co-actor is still a valid subtask. As a consequence, both possible complementary actions, A_1 and A_2 , in AEL are initially supported by excitatory input from the connected layers OML, IL and CSDL. As can be seen by the suprathreshold activation peak at position A_2 , the robot decides to first serve the human by grasping the wheel for handing it over (snapshots S4–S5). Later (see snapshots S6–S7) the robot selects and performs the action sequence associated with the achievement of its own subtask. The reasons for this ‘social’ attitude of the robot are slightly stronger synaptic connections from the intention layer IL to AEL compared to the synaptic links between representations in CSDL and AEL. As a consequence the decision process appears to be biased by the stronger input from the representation of the request. After both wheels have been attached to the platform, the human partner again demands an object (at time ≈ 78 s). This time the robot interprets the hand gesture as demanding a nut since the associated subtask representation in CSDL (‘fix wheel with nut’) has become active (see snapshot S8).



Fig. 8. Video snapshots that illustrate a scenario where the same observed action may have different interpretations, depending on the context in which it occurs, and thus lead to different action selection.

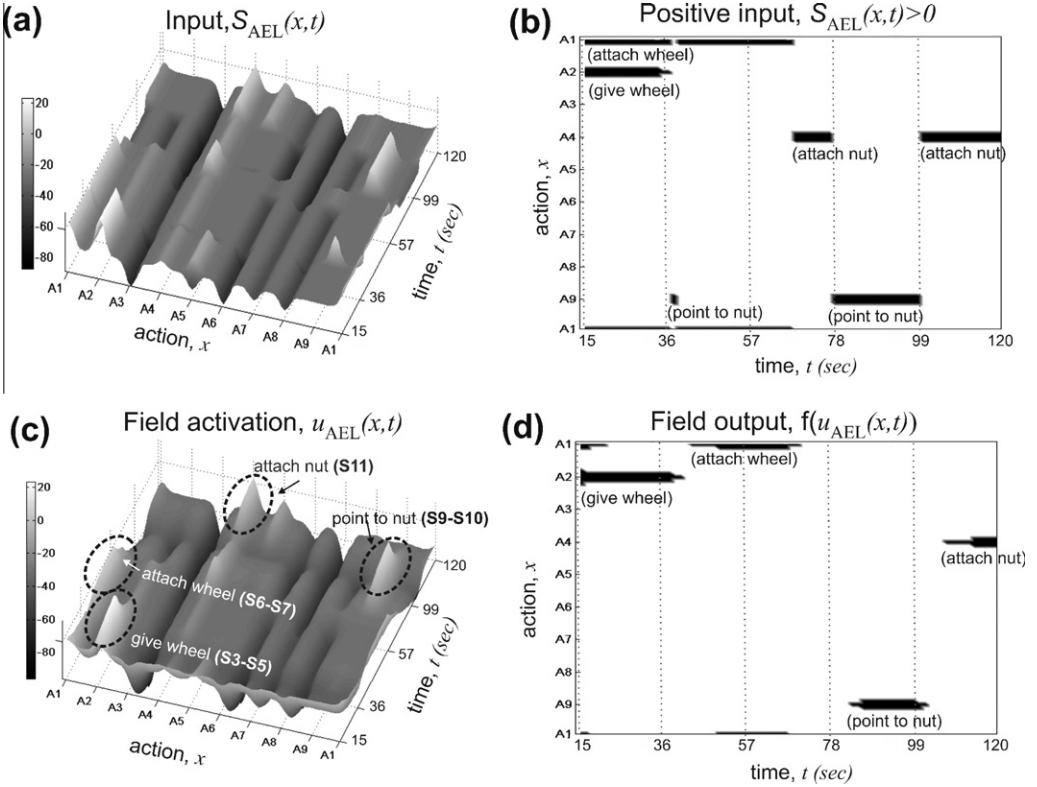


Fig. 9. Temporal evolution of total input to AEL (panels a and b) and field activation in time (panels c and d) for the experiment in Fig. 8.

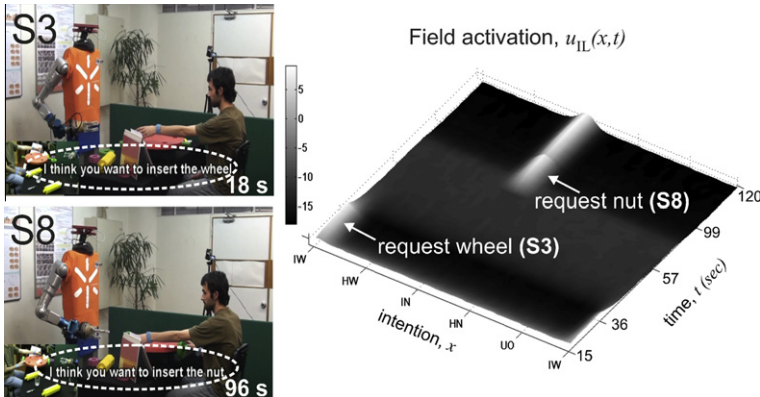


Fig. 10. Temporal activation in layer IL representing the inferred goal of the human partner for experiment in Fig. 8.

Fig. 10 illustrates the evolution of an activation peak in IL representing the inferred goal of the co-actor. However, a fulfilment of the request would not support efficient joint task execution since both co-actors have exactly one nut in their working areas. The information about the spatial distribution of the nuts represented in OML inhibits the representation of the complementary action that would

lead to the transfer of the nut (A_5) and triggers instead the selection of a pointing movement (A_9) as the most adequate motor behavior. Pointing towards the nut in the partner's workspace is an efficient way to attract attention and communicate the error to the co-actor (snapshots S9–S10). As illustrated by snapshots S11–S12 and the field activation after time 99 s, the robot subsequently decides to reach and grasp a nut with the purpose to fix the wheel at its construction side (A_4).

In this example pointing is an appropriate complementary behavior since the human partner could reach the nut that he had seemingly overlooked probably because of the presence of an obstacle. The situation would be different if the obstacle not only reduces the visibility of the nut but also makes it impossible for the user to grasp it. In this case, removing the obstacle or grasping the object to hold it out for the co-actor would be appropriate complementary actions (if the object could be reached by the robot). In the model, additional input from a population encoding the presence of an obstacle to the AEL could bias the decision process towards the selection of one of these behaviors. Interestingly, a class of mirror neurons has been recently described that are differently modulated by the location in space of observed motor acts relative to the monkey (Caggiano, Fogassi, Rizzolatti, Thier, & Casile, 2009). These are “grasping” neurons that become active when the experimenter places an object in the monkey's peripersonal space. In the experiments an obstacle was introduced that changed the properties of the mirror neurons according to the possibility that the monkey was able to interact with the object. The authors interpret their findings as further evidence for the hypothesis that mirror neurons encode observed actions for subsequent different types of behavior either direct grasping or intermediate steps like approaching the observed agent or removing the obstacle. This interpretation fits nicely to the highly context-sensitive mapping of observed actions onto executed actions implemented by the dynamic field model.

6.2. Anticipatory action selection

In the second experiment the robot has two nuts and one wheel in its workspace. The human partner has a wheel in his workspace and thus decides to start the assembly work by reaching and grasping the wheel on his side to directly attach it (Fig. 11). As can be seen in Fig. 11a, before the co-actor starts the movement the input to the decision field of the robot from the task representations in CGSL support the selection of the action sequence A_1 associated with a subgoal that the robot can achieve alone. However, immediately after motion onset the observed reaching behavior triggers a motor simulation process in ASL that anticipates that the co-actor most likely is going to attach the wheel on his side. Fig. 11a shows that a new input to the decision field appears in the period 10 and 14 s that is associated with the achievement of a future goal of the co-actor ('fix wheel with nut') represented in CSGL. Like in the preceding example the robot is supposed to serve first observed or anticipated needs of the human user. Since the input supporting a handing-over sequence (A_5) is stronger than the input supporting the goal-directed sequence A_1 , the robot decides at about time 18 s to reach and grasp a nut to hold it out for the co-actor (compare Fig. 11b and the snapshots in Fig. 11c). As discussed in Section 4 the predictive rather than reactive updating of common subgoals in CSGL is automatically triggered by the input from the representation of the inferred goal in IL.

6.3. Timing of actions matters

The last experiment highlights that the timing of actions is important for the coordination of decisions and intentions among the actors. Normally the faster actor takes the lead in the joint decision process and the slower actor follows by choosing actions that complement the observed ones. In the example shown in Fig. 12 the wheel on the co-actor's construction side has been already attached. The robot has two nuts in its workspace, whereas the co-actor has a wheel in his reach. The trial starts with the co-actor requesting an object (see snapshot S1). The robot infers that he wants to fix the wheel with a nut and decides to first hand over the nut (at time ≈ 14 s). Subsequently, the robot requests a wheel from the co-actor (at time ≈ 50 s). The first handing-over sequence thus appears to be initiated by the human, whereas the robot leads the transfer of the second object. The snapshots in Fig. 12c illustrate this sequence of human–robot interactions. Concerning the state of the construction and the initial distribution of object in the two working areas, the situation illustrated in Fig. 13c

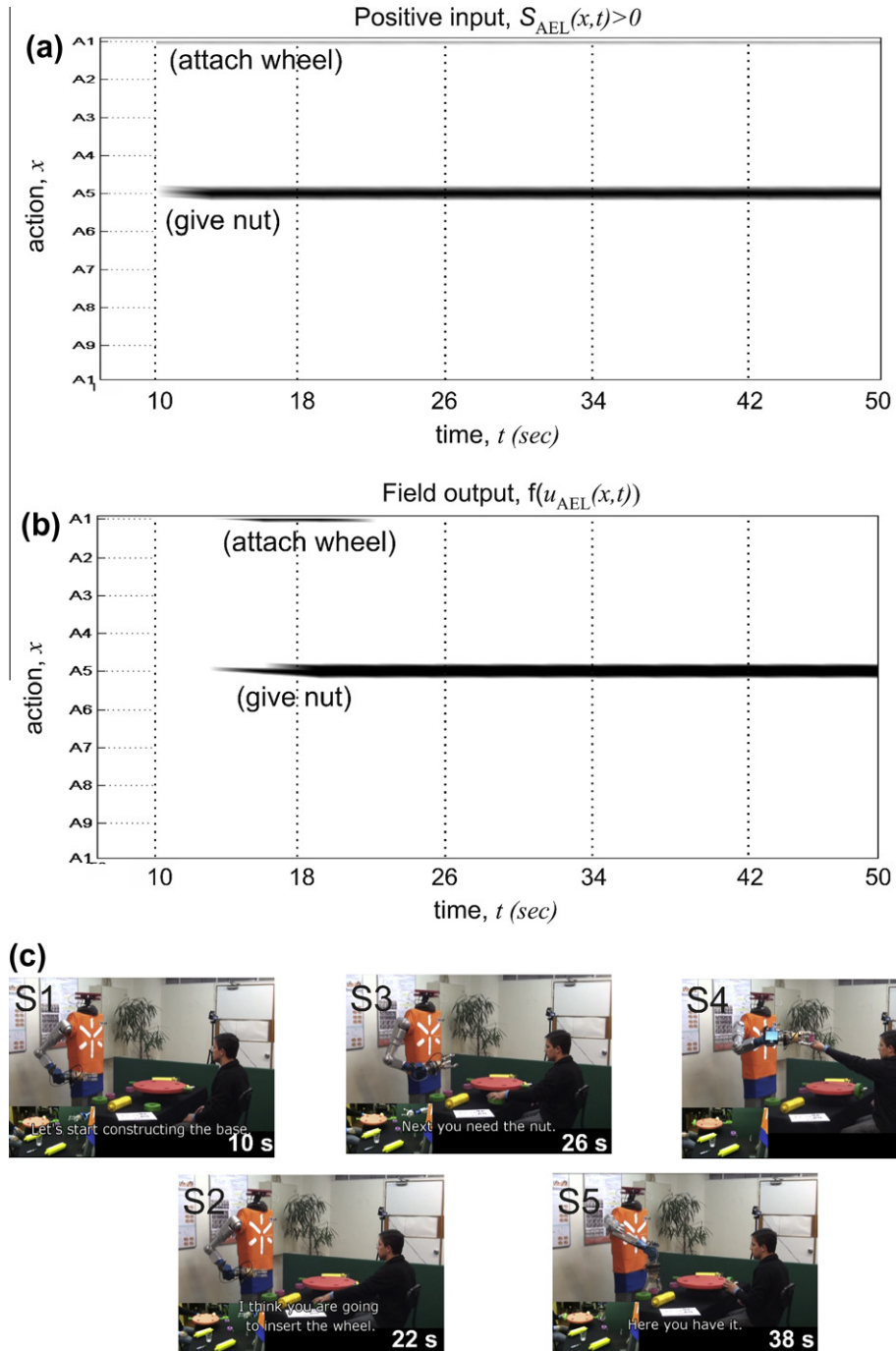


Fig. 11. Anticipatory action selection. (a) Top view of temporal evolution of total positive input to AEL showing the two possible actions, A_1 (reach and grasp wheel to attach) and A_5 (reach and grasp nut to handover). (b) Top view of temporal evolution of field output illustrating the competition between the two actions. Action A_5 wins the competition and the robot transfers the nut to the human partner. (c) Video snapshots (see text for details).

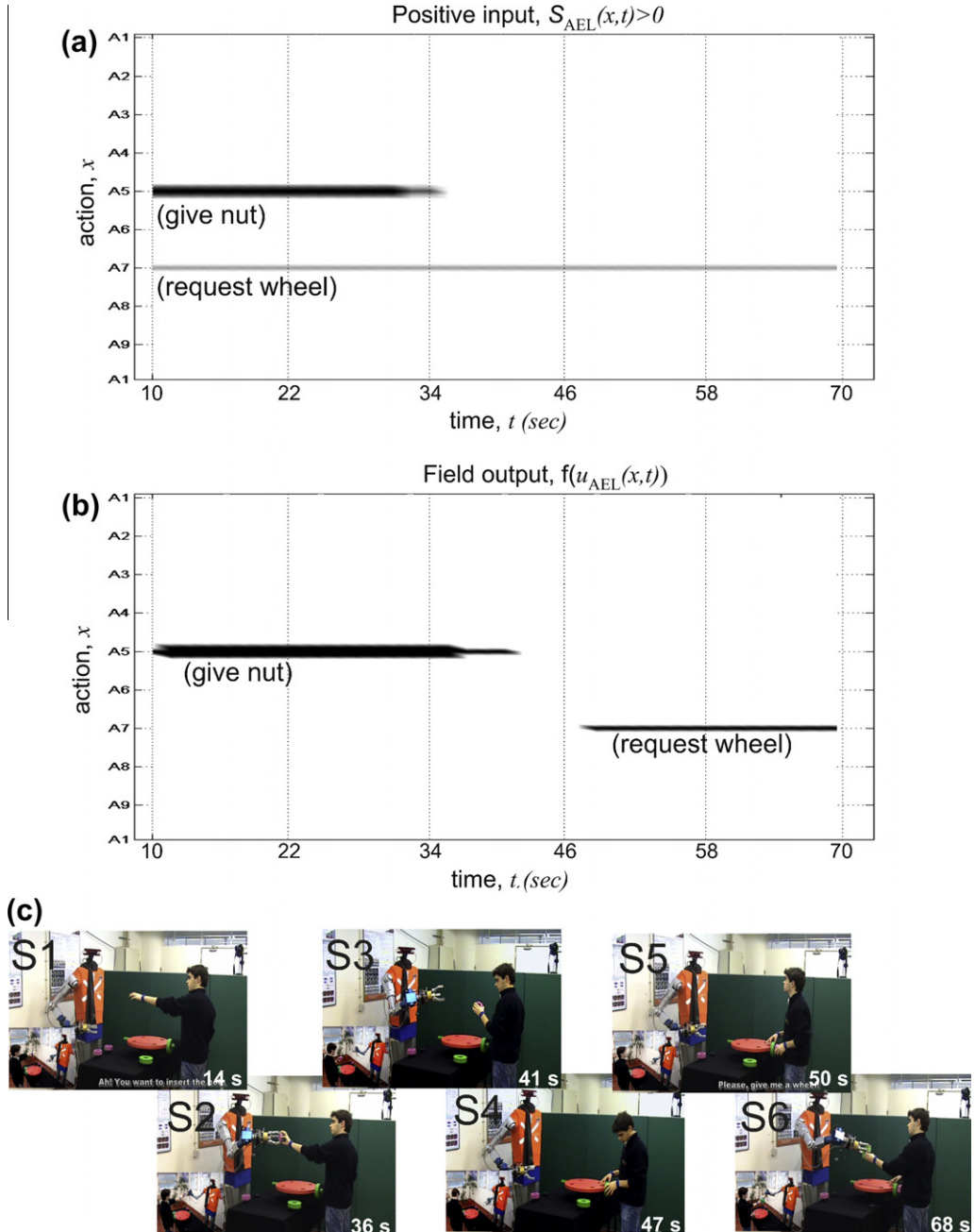


Fig. 12. Time matters. (a) Top view of temporal evolution of total positive input to AEL showing the two possible actions, A_5 (reach and grasp nut to handover) and A_7 (hold out hand as a request for a piece). (b) Top view of temporal evolution of field output illustrating the competition between the two actions. Action A_5 wins the competition and the robot gives the nut to the human partner and next action A_7 becomes activated and the robot request a wheel to the human. (c) Video snapshots (see text for details).

is exactly the same as in Fig. 12c. The only difference is that the co-actor now quickly tries to serve the needs of the robot after having fixed the wheel on his side with a nut. As shown in snapshots S4 and S5

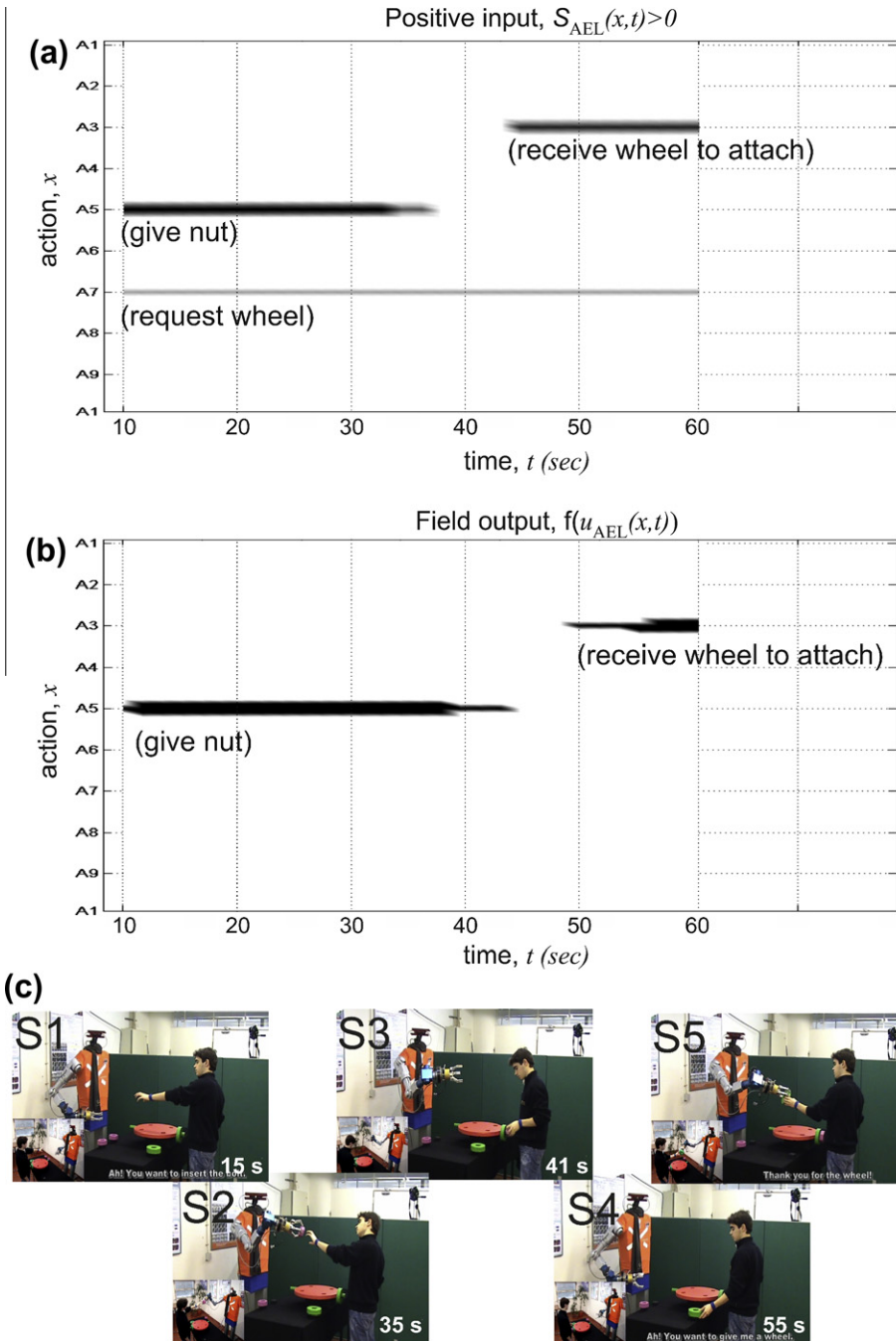


Fig. 13. Time matters. Situation similar to Fig. 12 but the human co-actor is now more aware of the needs of the robot and thus is faster to serve it. (a) Top view of temporal evolution of total positive input to AEL showing the possible actions, A_5 (reach and grasp nut to handover), A_7 (hold out hand as a request for a piece) and A_3 (receive and attach wheel) which appears because the robot infers that the partner wants to give a wheel. (b) Top view of temporal evolution of field output illustrating the competition between the actions and resulting decisions. (c) Video snapshots (see text for details).

the co-actor reaches and grasps a wheel with the intention to hold it out for the robot. Since through motor simulation the robot is able to predict at the time of the grasping the consequences of the ongoing action, it may prepare for receiving the wheel. The input associated with the respective action sequence appears at about 45 s (compare Fig. 12a). It is stronger compared to the input supporting the request of a wheel and thus biases the decision process of the robot.

7. Discussion

Decision making refers to the process of selecting a particular action from a set of alternatives. When acting alone an individual may choose a motor behavior that best serves a certain task based on the integration of sensory evidence and prior task knowledge. In a social context, this process is more complex since the outcome of one's decisions can be influenced by the decisions of others. A fundamental building block of social interaction is thus the capacity to predict and understand actions of others. This allows an individual to select and prepare an appropriate motor response which may range from cooperative to competitive (Sebanz et al., 2006).

Here we have presented a dynamic neural field model of decision making in a joint action task. In its multi-layered architecture the model reflects mechanisms that are believed to support the remarkably efficient and fluent interaction of humans in cooperative tasks. Most importantly, the model implements a highly context sensitive mapping of observed actions onto to-be-executed complementary actions. As our real-world joint construction experiments show, the robot responds to the same observed behavior in dependence of the context in which it occurs. This is in line with a growing body of experimental evidence supporting the notion that the matching of observed and executed actions in the mirror circuit is much more flexible than previously thought (Newman-Norlund, van Schie et al., 2007; van Schie et al., 2008). As the representation of context, goals and goal-directed action sequences are interconnected, the dynamic model explains how the observation of a motor act together with situational cues may directly activate the self-sustained population representations of the associated goal and the most appropriate complementary action. This 'automatic' process suggests that in known task settings the coordination of actions, goals and intentions in space and time between co-actors may occur rather effortlessly and does not require a fully developed human capacity for conscious control (Hassin, Aarts, & Ferguson, 2005; Ferguson & Bargh, 2004).

The theoretical framework of dynamic neural fields has been first introduced to the motor domain to model metric and dynamic aspects of motor planning and decision making found in neurophysiological and behavioral studies (for a recent review see Schöner (2008)). DNF-based models implement two basic ideas: (1) movement plans evolve continuously in time and are updated at any time during movement preparation as a function of new sensory evidence, and (2) the brain performs action selection and motor planning in an integrated manner (Gold & Shadlen, 2002). The activity in a neural population representing a certain decision variable increases continuously in time as a result of accumulated evidence represented by input from connected populations. If a certain activation threshold is reached the integration process is over and the system is committed to a decision. In the model this transition is paralleled by a transition from an input-driven to a self-stabilized regime of the field dynamics. The decision variable may represent simple movement parameters such as direction and extent (e.g., Erlhagen & Schöner (2002)) or like in the present application complete temporal behaviors such as grasping or entire action sequences composed of action primitives. The discovery of the mirror neurons in premotor cortex suggests that in motor planning areas neural populations encoding very different levels of abstraction coexist (Rizzolatti & Craighero, 2004).

Bayesian models represent a quite popular alternative approach for modeling decision and integration processes in the face of uncertainty (Kersten & Yuille, 2003; Körding & Wolpert, 2006). More recently, Bayesian inference and belief propagation have been used as theoretical tools to model also aspects of joint action coordination (e.g., Cuijpers, van Schie, Koppen, Erlhagen, & Bekkering, 2006; Hoffman & Breazeal, 2007). It is important to note that the dynamic field framework is compatible with central aspects of probabilistic models. For instance, the pre-activation below threshold of several populations in the action execution layer due to prior task knowledge and contextual information may be interpreted in the sense of a probability density function for different complementary actions.

This prior information has to be combined with evidence about the inferred goal of the co-actor. In fact, it can be shown that in the input-driven regime the field dynamics may implement Bayes' rules (Cuijpers & Erlhagen, 2008). In our view, there are two major advantages of the dynamic field approach. First, stabilizing decision against noise and fluctuations in the input stream is of particular importance in cases of high conflict between alternative complementary actions. Second, as an example of the dynamical approach to cognition (Schöner, 2008), a DNF-based model allows us to address the important temporal dimension of goal coordination in joint action (Sebanz et al., 2006). The decision process linked to complementary actions unfolds over time under multiple influences which are themselves modeled as dynamic representations with proper time scales. As our experiments show, the absence or delay of information from layer II will automatically lead to a decision that does not include the co-actor's behavior (or its interpretation). This may cause a change of roles in the joint task execution. Normally, the teammate with the faster decision process takes the lead in the cooperative task and the observer follows by choosing an action which complements the inferred goal (e.g., grasping the object with a complementary grip in the handing-over sequence). This flexibility in joint task execution greatly contributes to an efficient team performance.

Although the focus of the present study was on testing the dynamic model in real-time human–robot interactions it is worth mentioning that the model makes also predictions for human joint action coordination that could be further investigated in experiments. The reaction time study by van Schie and colleagues (2008) that showed evidence for an automatic facilitation of a complementary response in a cooperative setting used a one-to-one mapping between observed and to-be-executed actions. It would be interesting to extend this study to a more realistic situation in which a single observed action is compatible with several complementary behaviors like in the present assembly task. Due to the lateral inhibition in the action execution layer, the level of pre-activation of populations will decrease whenever several response alternatives are simultaneously supported by contextual and/or task information. Since the level of pre-activation affects the rate at which the population activity rises, the model predicts a dependence of reaction time on the number and probability of choices in the cooperative task (for a dynamic field approach to the classical Hick–Heyman law see Erlhagen and Schöner (2002)).

Among the many questions about joint action coordination that have not been directly addressed in the present experiments, an important one concerns how the model could be extended to deal with the coordination of multiple tasks. We have started to explore this challenge in a joint construction task in which the human–robot team has to assemble several distinct toy objects from a fixed set of components. Different from the present study, the robot does not directly participate in the construction work but serves as an intelligent assistant that pro-actively hands over components and informs the user about detected errors. By reducing the complexity of the action selection and execution process this choice allowed us to focus on the high-level planning aspects of the multi-object construction task. The dynamic field model implements the idea supported by many behavioral and neurophysiological studies that people encode goal-directed behaviors, such as assembling an object, by segmenting them into discrete actions, organized as goal-subgoal hierarchies (e.g., Hamilton & Grafton, 2008; Hard, Lozano, & Tversky, 2006). In the model, a population encoding a particular object as the final goal of the assembly work pre-activates through synaptic links the various populations encoding all associated subgoals. In the CGSL, the representations of the currently available subgoals for the team become activated above threshold due to the input from connected pools representing already achieved subtasks. Since different objects may share a set of subtasks (e.g., attaching two parts in a specific manner) synaptic links may exist to representations of subgoals belonging to objects other than the one currently under construction. These representations lack, however, the additional input from a population encoding the final goal. Consequently, their activation level remains below threshold. Using the same neural representations as part of action plans belonging to different tasks is attractive from an engineering point of view since it allows us to optimize the computational resources. It is important to stress however that joint action coordination may also benefit from special purpose representations (e.g., the grasping populations described in Fogassi et al. (2005)). The activation of intentional action chains in ASL during action observation supports the capacity of the robot to react in anticipation of the co-actor's motor intentions.

Even in routine tasks errors in joint action coordination may occur. The user may for instance select a part associated with a subgoal that has been already accomplished by the team or that represents an assembly step to be performed in the future only. It is thus important that the robot is able to cope with erroneous situations and unexpected events. In the present experiments, the robot points to a part in the co-actor's workspace that he has seemingly overlooked. The observed request gesture is thus associated not only with the complementary behavior of handing over the required part to the user but also with a hand movement that aims at attracting the co-actor's attention. The additional information about whether or not the part is located in the user's workspace (represented in the OML) biases the selection in the AEL. A more sophisticated action monitoring system that would allow the robot to deal with errors on the intention level should be able to detect a mismatch between predicted outcomes of observed actions and the possible subgoals for the team. Within the DNF-framework this can be achieved by postulating the existence of neural representations that integrate the activity from populations in the IL and the CSG. For instance, the co-actor's request of a part that is not compatible with the current state of the construction will automatically activate a population representing this mismatch. Through synaptic links to neural pools in the AEL, the suprathreshold activity of this population may in turn bias the selection of an adequate corrective response (Bicho, Louro, Hipolito, & Erlhagen, 2009).

For the present robot experiments all inter-field connections were tailored by the designer. Consequently, testing the dynamic field model of decision making in joint action was restricted to the specific assembly task. An important long-term goal of our research is to endow the robot with learning and adaptation capacities that ultimately will allow the artificial agent to autonomously develop the cognitive skills necessary for efficient joint action in new tasks from a minimal set of in-built representations. We adopt here a socially guided machine learning paradigm in which a human trainer teaches a robot through demonstration and verbal or gestural commands in much the same way as parents teach their children (e.g., Otero, Saunders, Dautenhahn, & Nehaniv, 2008; Thomaz & Breazeal, 2008). First experimental results of our attempt to apply a learning dynamics for establishing inter-field connections show the feasibility of the approach. Using correlation based learning rules with a gating that signals the success of behavior, we have shown for instance how mirror-like representations that support an action understanding capacity may develop during learning and practice (Erlhagen et al., 2006a; Erlhagen et al., 2006b). Importantly, the developmental process goes beyond a simple modification of parameters of pre-defined representations. It may explain the emergence of new task-specific populations which have not been introduced to the architecture by the human designer (Erlhagen et al., 2007).

A major goal of our group is to advance towards a new generation of robots able to interact with humans in a more natural and efficient manner. We believe that taking inspiration from biology to make both the observable trajectories and the cognitive processes supporting joint action more human-like is a promising approach since it makes the artificial agent more predictable for the human user. While many technical aspects of robotics (e.g., vision, sensorimotor coordination) have been simplified, we also believe that the robotics work based on the unifying framework of dynamic neural fields is potentially very interesting for researchers from the cognitive and neuroscience domains. Implementing the dynamics of cooperative joint action in an embodied cognitive system allows them to directly test their theories and hypothesis about joint action coordination.

Acknowledgments

The present research was conducted in the context of the fp6-IST2 EU-IP Project JAST (Project No. 003747) and partly financed by the FCT grants POCI/V.5/A0119/2005 and CONC-REEQ/17/2001. We would like to thank Profs. Harold Bekkering and Ruud Meulenbroek for the numerous discussions, and Emanuel Sousa, Flora Ferreira, Nzoji Hipolito, Rui Silva and Toni Machado for their help during the robotic experiments. We also thank the anonymous reviewers for their insightful comments and feedback on how to improve the manuscript.

Appendix A. Supplementary data

Supplementary data associated with this article, like numerical values for the Joint Action Model parameters, can be found in the supplemental material.

Supplementary data associated with this article can be found, in the online version, at doi:10.1016/j.humov.2010.08.012.

References

- Amari, S. (1977). Dynamics of pattern formation in lateral-inhibitory type neural fields. *Biological Cybernetics*, 27, 77–87.
- Bicho, E., Louro, L., & Erlhagen, W. (2010). Integrating verbal and nonverbal communication in a dynamic neural field architecture for human–robot interaction. *Frontiers in Neurobotics*, 4. doi:10.3389/fnbot.2010.00005.
- Bicho, E., Louro, L., Hipolito, N., & Erlhagen, W. (2008). A dynamic neural field architecture for flexible and fluent human–robot interaction. In *Proceedings of the 2008 international conference on cognitive systems* (pp. 179–185). Germany: University of Karlsruhe.
- Bicho, E., Louro, L., Hipolito, N., & Erlhagen, W. (2009). A dynamic field approach to goal inference and error monitoring for human–robot interaction. In K. Dautenhahn (Ed.), *Proceedings of the 2009 international symposium on new frontiers in human–robot interaction. AISB 2009 Convention* (pp.31–37) Edinburgh; Heriot-Watt University.
- Bicho, E., Mallet, P., & Schöner, G. (2000). Target representation on an autonomous vehicle with low-level sensors. *The International Journal of Robotics Research*, 19, 424–447.
- Blakemore, S. J., & Decety, J. (2001). From the perception of action to the understanding of intention. *Nature Reviews Neuroscience*, 2, 561–567.
- Brass, M., Bekkering, H., & Prinz, W. (2001). Movement observation affects movement execution in a simple response task. *Acta Psychologica*, 106, 3–22.
- Caggiano, V., Fogassi, L., Rizzolatti, G., Thier, P., & Casile, A. (2009). Mirror neurons differentially encode the peripersonal and the extrapersonal space in monkey. *Science*, 324, 403–406.
- Cuijpers, R., & Erlhagen, W. (2008). Implementing Bayes' rules with neural fields. In V. Kurkov (Ed.), *ICANN 2008, part II. LNCS 5164* (pp. 228–237). Springer Verlag.
- Cuijpers, R. H., van Schie, H. T., Koppen, M., Erlhagen, W., & Bekkering, H. (2006). Goals and means in action observation: A computational approach. *Neural Networks*, 19, 311–322.
- di Pellegrino, G., Fadiga, L., Fogassi, L., Gallese, V., & Rizzolatti, G. (1992). Understanding motor events: A neurophysiological study. *Experimental Brain Research*, 91, 176–180.
- Dominey, P. F., & Warneken, F. (in press). The basis of shared intention in human and robot cognition. *New Ideas in Psychology*. doi:10.1016/j.newideapsych.2009.07.006.
- Erlhagen, W., & Bicho, E. (2006). The dynamic neural field approach to cognitive robotics. *Journal of Neural Engineering*, 3, R36–R54.
- Erlhagen, W., Mukovskiy, A., & Bicho, E. (2006a). A dynamic model for action understanding and goal-directed imitation. *Brain Research*, 1083, 174–188.
- Erlhagen, W., Mukovskiy, A., Bicho, E., Panin, G., Kiss, C., Knoll, A., et al. (2006b). Goal-directed imitation for robots: A bio-inspired approach to action understanding and skill learning. *Robotics and Autonomous Systems*, 54, 353–360.
- Erlhagen, W., Mukovskiy, A., Chersi, F., & Bicho, E. (2007). On the development of intention understanding for joint action tasks. In *6th IEEE international conference on development and learning* (pp. 140–145). Imperial College London.
- Erlhagen, W., & Schöner, G. (2002). Dynamic field theory of movement preparation. *Psychological Review*, 109, 545–572.
- Ferguson, M. J., & Bargh, J. A. (2004). How social perception can automatically influence behavior. *Trends in Cognitive Sciences*, 8, 33–39.
- Flash, T., & Hogan, N. (1985). The coordination of arm movements: An experimentally confirmed mathematical model. *Journal of Neuroscience*, 5, 1688–1703.
- Fogassi, L., Ferrari, P. F., Gesierich, B., Rozzi, S., Chersi, F., & Rizzolatti, G. (2005). Parietal lobe: From action organization to intention understanding. *Science*, 308, 662–667.
- Fong, T., Nourbakhsh, I., & Dautenhahn, K. (2003). A survey of socially interactive robots. *Robotics and Autonomous Systems*, 42, 143–166.
- Gold, J. I., & Shadlen, M. N. (2002). Banburismus and the brain: The relationship between sensory stimuli, decisions, and reward. *Neuron*, 36, 299–308.
- Grossberg, S. (1973). Contour enhancement, short term memory and constancies in reverberating neural networks. *Studies in Applied Mathematics*, 52, 217–257.
- Hamilton, A. F., & Grafton, S. T. (2008). Action outcomes are represented in human inferior frontoparietal cortex. *Cerebral Cortex*, 18, 1160–1168.
- Hard, B. M., Lozano, S. C., & Tversky, B. (2006). Hierarchical encoding of behavior: Translating perception into action. *Journal of Experimental Psychology: General*, 135, 588–608.
- Hassin, R. R., Aarts, H., & Ferguson, M. J. (2005). Automatic goal inferences. *Journal of Experimental Social Psychology*, 41, 129–140.
- Hoffman, G., & Breazeal, C. (2007). Cost-based anticipatory action selection for human–robot fluency. *IEEE Transactions on Robotics*, 23, 952–961.
- Kersten, D., & Yuille, A. (2003). Bayesian models of object perception. *Current Opinion in Neurobiology*, 13, 1–9.
- Körding, K. P., & Wolpert, D. M. (2006). Bayesian decision theory in sensorimotor control. *Trends in Cognitive Sciences*, 10, 319–326.
- Laing, C. R., Tray, W. C., Gutkin, B., & Ermentrout, G. B. (2002). Multiple bumps in a neural model of working memory. *SIAM Journal of Applied Mathematics*, 63, 62–97.
- Lommertzen, J., Costa e Silva, E., Cuijpers, R., & Meulenbroek, R. (2008). Collision-avoidance characteristics of grasping: Early signs in hand and arm kinematics. In *Anticipatory behavior in adaptive learning systems (ABIALS 2008) post-conference proceedings, Munich, Germany* (pp. 188–208).
- Meulenbroek, R., Rosenbaum, D., Jansen, C., Vaughan, J., & Vognt, S. (2001). Multijoint grasping movements: Simulated and observed effects of object location, object size and initial aperture. *Experimental Brain Research*, 2, 219–234.

- Newman-Norlund, R. D., Noordzij, M. L., Meulenbroek, R. G. J., & Bekkering, H. (2007). Exploring the basis of joint action: Coordination of actions, goals and intentions. *Social Neuroscience*, 2, 48–65.
- Newman-Norlund, R. D., van Schie, H. T., van Zuijlen, A. M. J., & Bekkering, H. (2007). The mirror neuron system is more active during complementary compared with imitative action. *Nature Neuroscience*, 10, 817–818.
- Otero, N., Saunders, J., Dautenhahn, K., & Nehaniv, C. L. (2008). Teaching robot companions: The role of scaffolding and event structuring. *Connection Science*, 20, 111–134.
- Rizzolatti, G., & Craighero, L. (2004). The mirror–neuron system. *Annual Review of Neuroscience*, 27, 169–192.
- Rizzolatti, G., Fogassi, L., & Gallese, V. (2001). Neurophysiological mechanisms underlying the understanding and imitation of action. *Nature Reviews*, 2, 661–670.
- Rosenbaum, D., Meulenbroek, R., Vaughan, J., & Jansen, C. (2001). Posture-based motion planning: Applications to grasping. *Psychological Review*, 108, 709–734.
- Schöner, G. (2008). Dynamical systems approaches to cognition. In R. Sun (Ed.), *The Cambridge handbook of computational psychology* (pp. 101–125). Cambridge University Press.
- Schöner, G., Kopecz, K., & Erlhagen, W. (1997). The dynamic neural field theory of motor programming: Arm and eye movements. In P. Morasso & V. Sanguineti (Eds.), *Self-organization, computational maps and motor control* (pp. 271–310). Elsevier Science.
- Sebanz, N., Bekkering, H., & Knoblich, G. (2006). Joint action: Bodies and minds moving together. *Trends in Cognitive Sciences*, 10, 70–76.
- Silva, R., Bicho, E., & Erlhagen, W. (2008). Aros: An anthropomorphic robot for human–robot interaction and coordination studies. In *8th Portuguese conference on automatic control (Controlo 2008)* (pp. 819–826). 21–23 July, Vila Real, Portugal.
- Thomaz, A. I., & Breazeal, C. (2008). Teachable robots: Understanding human teaching behaviour to build more effective robot learners. *Artificial Intelligence*, 172, 716–737.
- van Schie, H. T., van Waterschoot, B. M., & Bekkering, H. (2008). Understanding action beyond imitation: Reversed compatibility effects of action observation in imitation and joint action. *Journal of Experimental Psychology: Human Perception and Performance*, 34, 1493–1500.
- Vaughan, J., Rosenbaum, D., & Meulenbroek, R. (2006). Modeling reaching and manipulation in 2- and 3-d workspaces: The posture-based model. In *Proceedings of ICDL 2006* (pp. cd-rom). June 3, Bloomington.
- Westphal, G., von der Malsburg, C., & Würtz, R. P. (2008). Feature-driven emergence of model graphs for object recognition and categorization. In A. Kandel, H. Bunk, & M. Last (Eds.), *Applied pattern recognition* (pp. 155–199). Springer Verlag.
- Wilmzig, C., Schneider, S., & Schöner, G. (2006). The time course of saccadic decision making: Dynamic field theory. *Neural Networks*, 19, 1059–1074.
- Wilson, H. R., & Cowan, J. D. (1973). A mathematical theory of the functional dynamics of cortical and thalamic nervous tissue. *Kybernetik*, 13, 55–80.
- Wilson, M., & Knoblich, G. (2005). The case for motor involvement in perceiving conspecifics. *Psychological Bulletin*, 131, 460–473.