

From egocentric systems to systems allowing for Theory of Mind and mutualism

Holk Cruse¹ and Malte Schilling^{1,2}

¹Biological Cybernetics and Theoretical Biology Department,
and Center for Cognitive Interaction Technology (CITEC), University of Bielefeld, Germany

²ICSI Berkeley
holk.cruse@uni-bielefeld.de

Abstract

Simple artificial agents representing more or less elaborated Braitenberg vehicles, usually adopt an egocentric view. One example is Walknet, a biologically inspired neural network controlling hexapod walking. Here we show how such a controller can be expanded to be able to interpret observed behaviours that are performed by other individuals, i.e. the system shows properties of a mirror system. This allows to further expand the network to become an “allocentric” system that might implement subjective feelings which could be attributed to other individuals, i.e. the system implements a Theory of Mind. As a last expansion we introduce a two-body model, or we-model, which may allow for mutualism. Application of we-models allows for what often has been called the third person’s view. The different steps proposed can be interpreted as corresponding to an evolutionary development.

Introduction

Artificial agents being based on natural creatures may usually be characterized as to hold an ‘egocentric’ view: in such agents, the sensory input is related to the own body representing the center of the agent’s world. Correspondingly, motor output activities are based on the own geometrical—and possibly mental—position. Here we attempt to introduce a way how the controller of such an autonomous agent may be changed to allow the agent to ‘put itself into the partner’s shoes’, in other words to allow for theory of mind (ToM), and to show empathy. A further goal is to develop a (neuronal) control structure that may form the basis of mutualism, i.e. the faculty to cooperate with a partner using shared goals (Tomasello, 2009). Such a control structure may serve as a quantitatively defined hypothesis and may as such help to understand the underlying mechanisms of the corresponding biological system.

When attempting to simulate higher mental functions as are specific memory systems, attention, cognition or consciousness, for example, authors do, in general, not apply a whole-systems approach, but instead consider specific networks suited to represent the specific function of interest. Therefore, in many cases, it remains open how these specific networks may be embedded into the complete system, i.e. how the different networks are switched on or off and how these local networks receive input from and provide output for the complete system. To avoid this problem, we take a whole-systems approach. We investigate such phenomena under the

condition that these networks are embedded into an autonomously behaving agent, i.e. an agent equipped with a body characterised by many parallel and serially arranged degrees of freedom and a control network containing a set of preexisting reactive behaviours.

Schilling and Cruse (submitted) have proposed a network that has been worked out in more detail and called reaCog (this work is based on the reactive control system Walknet (Dürr et al., 2004) and the cognitive extensions have been introduced in Cruse and Schilling (2010)). This network is able to control a hexapod system by applying a structure consisting of two levels. The lower level is endowed with properties that correspond to insect-type behaviours (as are walking, climbing and navigation), about which already detailed knowledge is available (Dürr et al. 2004, Bläsing 2006, Wehner 2008). This level is based on a reactive, or behaviour-based, architecture, i.e. a collection of local, in general recurrent, neural networks (RNN). The second level of reaCog concerns an expansion allowing for the introduction of cognitive abilities as explained below. Generally, the architecture of our system is not based on the idea to consist of one holistic RNN, but represents a localist approach the advantages of which are convincingly advocated by Cooper and Shallice (2006).

When starting with an insect-like body and insect-inspired behaviour-based networks we do not imply that insects were endowed with higher cognitive functions as are metacognition, ToM or consciousness, although already in insects a number of astonishing properties can be found which by some authors are called cognitive (e.g. application of concepts like symmetry, sameness or protocounting, see Menzel et al., 2007). However, we assume that any cognitive system is strongly relying on such reactive—or behaviour-based—structures. Different to a reactive system, a cognitive system in the strict sense should be able to exploit stored information independent of the context in which this information has been acquired. This means, a cognitive system should be able to combine existing memory elements in a new way and use these new combinations for controlling behaviour and planning ahead. As we have shown by having developed reaCog, only a limited number of expansions are required to reach such a cognitive level (Cruse and Schilling 2010; Schilling and Cruse, submitted). The most important expansion concerns the introduction of a ‘manipulable’ body model. In order to be able to plan ahead, this internal model of

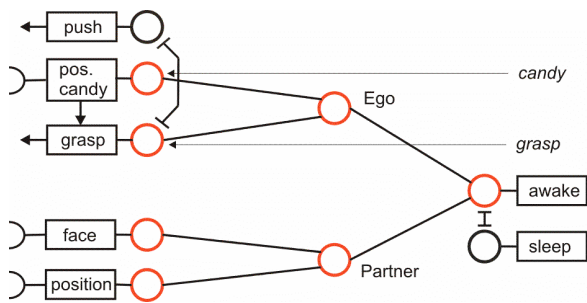


Fig. 1. An egocentric network represents the situation “Ego grasp candy”. The figure shows a section of the network *reaCog* (Schilling and Cruse, *subm.*). Local networks are symbolized by rectangles and names. Motivation units are shown by circles (connection to the corresponding network see Fig. 4). Active motivation units are marked by red colour. Arrows represent excitatory connections, T-shaped connections are inhibitory. Visual and proprioceptive input is marked by the half-circles, left side. Acoustic input representing words is shown by italic letters at the right side.

the own body (plus some aspects of the world, e.g. an obstacle) is required to internally simulate different behaviours in order to test whether this specific behaviour is suited to cope with an actual problem. The second expansion concerns an attention system. This system consists of two layers, a spreading activation layer (SAL) and a winner-take-all layer (WTA). This two-layer network enables the agent to select a specific behavioural element, which is normally not activated in the actual context. Via internal simulation, the system can then test whether this newly selected behavioural element is suited to solve the problem at hand, a procedure that has been termed “*probehandeln*” following Freud (1911). New behaviours found by this procedure and that, by means of the simulation and the subsequent behavioural test, prove to be adaptive will be stored in the long-term memory, thereby enriching the behaviour-based architecture. As for a well designed reactive system new problems may occur only rarely, *reaCog* can be regarded a reactive system that exploits its cognitive properties only for short periods of time required to solve a problem at hand.

Based on the ideas of Narayanan (Narayanan, 1997 and Feldman and Narayanan, 2004) and Steels (1995, 2003) we have further designed a simple expansion of *reaCog* that allows connecting behavioural elements of this system with so called word nets (RNNs representing an individual verbal expression, e.g. “leg”, or “swing”) that carry the corresponding meaning (Cruse, 2010). Therefore, the symbols are grounded (Steels, 2003) allowing the agent to ‘understand’ the meaning of such a word when given to the agent.

Like most other autonomous systems, *reaCog* holds an “egocentric” view. The agent might be able to recognize and represent objects. We further assume that the agent can also recognize, as a specific kind of object, a conspecific (see Steels and Spranger, 2008 and Spranger et al., 2009 for solutions). In addition we assume that the agent can attribute properties to the object or the partner (e.g. a face, a spatial position). All these expansions, however, do not enable the agent to “put himself into the partner’s shoes”. In other words,

the agent is not able to realize that the partner may see the agent himself as having a property (e.g. a position). Thus, in this network there is no possibility to represent the change of roles (“If I were him”). In other words, the capability to have a ToM is lacking. A classical procedure for testing whether an agent allows for the ability of ToM is the so called Sally-Anne task. Two subjects are shown that a candy lying on the table is hidden under a black cover. Then one subject, Sally, has to leave the room whilst the candy is now hidden under the white cover, as observed by Anne. After Sally has come back, Anne is asked under which cover Sally will probably search for the candy. If Anne points to the black cover, she is assumed to have ToM, but not, if she points to the white cover where the candy really is placed.

The network *reaCog* even less shows the ability to perform mutualistic behaviour (Tomasello, 2009), i.e. to develop shared goals and to try to follow them, even when the individual agent may receive no specific advantage. A simple example is when two individuals are trying to carry a load, for instance a table through an environment containing obstacles. In the reminder we show how *reaCog* can be expanded to endow the agent with these capabilities. To be in a position to explain the structures and their properties in an easily understandable way, we illustrate the expansions of *reaCog* by attempting to maintain the number of neuronal units as small as possible. In this way we hope to provide a functional understanding of how systems able to develop a ToM and later a structure allowing for mutualism may have arisen from an egocentric system. The different steps introduced might represent a hypothetical evolutionary sequence.

The Model

To simplify the description, we will focus on a small section of *reaCog* as illustrated in Fig. 1. Basically, the network consists of sensorimotor networks, or memory elements, connected with motivation units. In the figures, the networks are indicated by rectangles with verbal descriptors. Motivation units (depicted as circles) can adopt an activation value within the interval $[0,1]$. In the figures, activated units are marked as red circles, inactive ones are shown as black circles. Two of these motivation units may either be connected via (mutual) inhibition or via (mutual) excitation, or not be connected at all. Groups of excitatorily connected units stabilize each other. I.e. when one unit of such a group is activated, all the members of that group will become activated, too, except for those units that are connected via mutual inhibition. These inhibitory connections form a local winner-take-all (WTA) net with the consequence that only one of these units will stay active over some iterations. Two such motivation units represent the state Awake and the state Sleep, respectively. In the awake state, several sensory or motor elements are activated. These elements may form different contextual groups. Here we focus on two such groups. One group refers to external objects, in this case a conspecific (“partner”), represented by the memory elements “face” and “position”, which stand for the visual appearance and spatial location of the partner to be recognized by the corresponding networks. Together with the unit Partner these motivation units form an excitatory network. The elements of the second group refer to the agent. The agent can select between a number of actions

(in Fig. 1 “push” and “grasp”), the motivation units of which are connected via mutual inhibition (connections with T-shaped endings). The agent is also assumed to recognize an object, a candy lying on the table. Fig. 1 shows a memory element representing the position of the candy (*pos.candy*) relative to the agent. The agent may also be equipped with a network representing the experience of pain, which is connected to any specific body position, but this faculty will only be explained later. The motivation unit connecting the agent-related elements has been called Ego unit in the figures. To avoid a possible misunderstanding, it should be made clear that this name represents only a technical term and should not be understood as to mean that the agent has any kind of self-knowledge. As mentioned, the system may also be equipped with word nets that allow to recognize verbal statements as “*grasp*” or “*candy*” or “*partner*” which, if stimulated, activate the corresponding sensorimotor networks (in the figures these inputs are indicated by the terms given in italic; the word nets themselves are not shown).

Of course, any partner, if being equipped with a corresponding network, may likewise recognize our agent, but, as mentioned, the agent does not know this.

The behaviour-based—or sensorimotor—RNNs indicated by rectangles in the figures might be realized as simple associators connecting a sensory input with a motor output (Dürr et al., 2004; Cruse and Wehner, 2011) and may function as an implicit body model, that can be used to control the behaviour by computing the inverse kinematics. Alternatively, as conceptualized in *reaCog* (Cruse and Schilling, 2010; Schilling, 2011; Schilling and Cruse, *subm.*), sensorimotor RNNs may be connected to an explicit body model. In this case, the network is equipped with a switch that allows to turn on or off the motor output to either control the behaviour or instead to activate only the body model and in this way simulate the behaviour. In the latter case, the system may be termed to imagine this behaviour.

To realize the motivation units and RNN units we use the so called Input Compensation (IC) units, type suppression units (Kühn et al. 2007, Makarov et al. 2008), first, because a simple learning algorithm is available to train such networks. Secondly, because such networks maintain the input activation as long as the input is provided, but, if trained to hold a static attractor, also after the input is switched off. A motivation unit that is connected to a behaviour-based RNN, controls the output of its network by multiplying the output by its activation value (see below, Fig. 5). In this way, a motivation unit when activated may be called to ‘open’ the corresponding network (representing a top-down influence). As will be mentioned below, sensorimotor networks may also be used to respond to sensory input. In this case, the network showing the best fit to the actual sensory input (or the smallest error) will activate its motivation unit (this bottom-up influence is not depicted in Fig. 5). In the simulation proposed here, only the motivation unit network has been studied (for an explicit simulation of such a network see Cruse and Wehner, 2011).

Phenomenal aspect: Before we continue to describe the property of the network in more detail, a fundamental, and unsolved problem has to be addressed. When trying to understand a cognitive system the question arises how a neuronal system representing a physical structure is able to

allow for the faculty to experience subjective feelings, an example is feeling pain. This subjective or phenomenal aspect is relevant for (at least some) living systems. What is the problem? We can easily think of neuronal structures that, activated by nociceptors, for example, may produce chemical substances or activate specific behaviours (e.g. withdrawal or speech acts), i.e. form a series of causally connected physical states. But there is no concrete idea how (and why) the fact that these (or some of these) physical activities are accompanied by the feeling of pain, i.e. the subjective aspect, may be reified. The problem of understanding the relation between the physical aspect and the phenomenal aspect has eventually been termed the ‘hard problem’ (Chalmers, 1996) and will not attempted to be solved here. In order to be nevertheless able to use terms describing (or at least associated with) subjective feelings when discussing the properties of our network, we make the following assumption. An RNN as used here can adopt attractor states that are reached when the network has been given enough time for relaxation. In mathematical terms the attractor state can be defined as the so-called harmony value of the net reaching a maximum value. Following Cruse (1999, 2003) we assume that the activation of such a network is accompanied by subjective experience (or a phenomenal aspect) if the harmony value of the net has reached a given threshold, in other words, if the net has sufficiently well approached its attractor state. This hypothesis does of course not represent a solution of the hard problem, but nonetheless provides a way to operationalise the problem. Its function in this context is to allow us using terms associated with subjective or phenomenal aspects when describing states of our physical network. Using this hypothesis we are in a position to bridge the ‘explanatory gap’ on a descriptive level. If other mechanisms underlying the phenomenal aspect were found, they could replace our hypothesis without, as we believe, influencing the rest of the arguments.

The functioning of the network – an example

The agent equipped with *reaCog*, the, for our discussion, relevant part of which is depicted in Fig. 1, is able to show the following simple behaviour. If we assume that elements “*grasp*” and “*pos.candy*” are activated by an external verbal command as indicated by thin arrows (in the figures marked by italic letters, e.g. Fig. 1 *grasp*, *candy*), this input will activate the motivation units *grasp* and *pos.candy*. The former will open the behaviour represented in the RNN *grasp* and activate the unit *Ego*. Further, the unit *pos.candy* when activated will open the RNN allowing to recognize the spatial position of the candy. The *grasp* network receives input from the *pos.candy* network that provides the information to the *grasp* network concerning the goal for the movement to be performed. Therefore, the movement can now be executed. As an alternative to verbally given input, the agent, after having registered the candy, may decide to perform a grasp movement, the decision being determined by its internal state requiring a network not shown in the figures. In the following examples we will however use verbal input only, because this simplifies explanation of the concepts proposed. As illustrated in Fig. 1, at the same time the agent may be able to represent a partner, characterized by its face and its position.

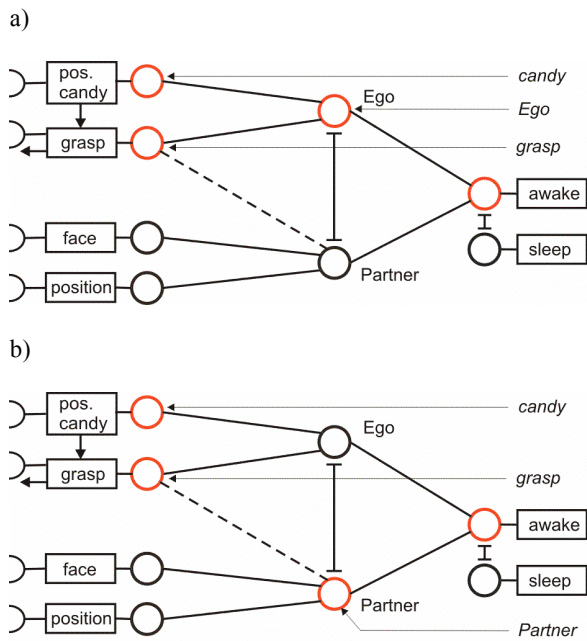


Fig. 2. An egocentric network representing the situation “Ego grasp candy” (a) and the situation “Partner is seen as grasping a candy” (b). The sensorimotor element “grasp” provides motor output and receives sensory (e.g. visual) input. Its units show properties corresponding to those of mirror neurons as it represents a circuit shared between the Ego and the partner units. See Fig. 1 for further explanation.

Mirror systems

How may this network be changed to allow for ToM and mutualism? Several changes are proposed as will be illustrated in consecutive steps depicted in Figs. 2, 3 and 4. A body model, apart from being used to control movement by calculating the inverse kinematics (Fig. 1), can also be used for a different purpose. When observing somebody else performing a grasp or a push movement, the visual input can be given to the body model which then can be used to simulate, or “internally copy”, the observed behaviour (e.g. “grasp”) following the “simulation theory” (e.g. Jeannerod, 2006 & 1999, Gallese & Lakoff, 2005). This application of the body model is suited to minimize errors when interpreting the (underspecified) visual input (e.g. Schilling, 2011). To symbolize this ability, in Fig. 2 the net ‘grasp’ is also equipped with sensory (visual) input. By application of a specific RNN forming a holistic system as has been proposed by Cruse and Schilling (2010) and Schilling (2011), one and the same body model is exploited for both purposes as are motor control and interpretation of sensory input. If a grasping movement is observed, the body model activates the element ‘grasp’. To allow the representation of the partner performing a grasping movement, too, we need another expansion, namely the introduction of connections between the unit representing the partner with (some of) the behavioural elements that, in the egocentric system (Fig. 1), are only

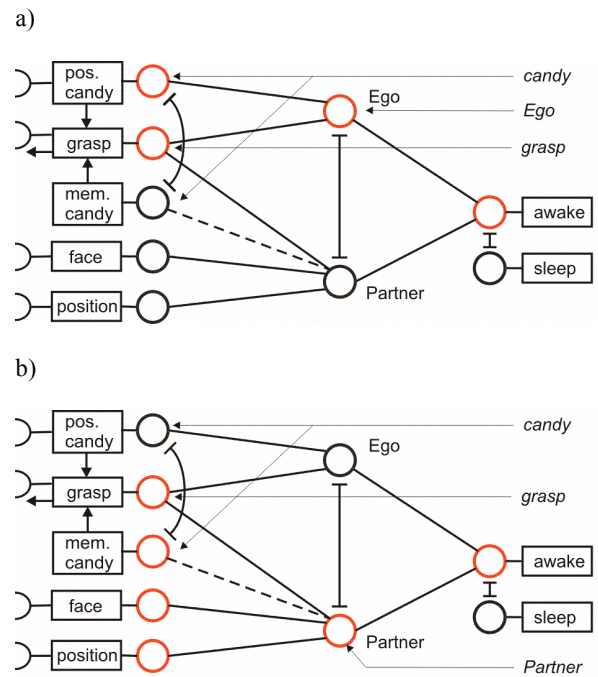


Fig. 3. A network being able to represent an egocentric view (a, situation “Ego grasp candy”) and the view as seen by the partner (b, situation “Partner grasp candy”), thus allowing for ToM. For further explanations see Fig: 1 and text.

connected with the Ego Unit. In our example this refers to element ‘grasp’ (see Fig. 2, dashed line). In addition, Unit Ego and unit Partner have to be connected via mutual inhibition (Fig. 2). This means that either unit Ego or unit Partner can be activated at a given moment in time.

With this network we can represent two situations: (i) if, as depicted in Fig. 2a, units Ego, grasp and pos.candy are coactivated, the network represents the agent to grasp the candy or to imagine such a grasping movement (the representation of this situation is already possible for the network shown in Fig. 1). (ii) However, the agent can also record a grasping movement of the partner. In this case, the sensorimotor element ‘grasp’ is activated together with the unit Partner, whereas unit Ego is inhibited. In Fig. 2b this situation is illustrated by motivation unit Partner shown in red and unit Ego in black. In both situations the neurons of the element grasp are activated. Such an architecture has eventually be termed to apply ‘shared circuits’ and strongly reminds of properties characterizing mirror neurons. Therefore, application of such shared circuits has been described as ‘mirroring’ (Keysers and Gazzola, 2011). Units of the grasp net represent to movement and its goal, and thus correspond to represent a motor act as attributed to mirror neurons (Rizzolatti and Luppino 2001). However, the goal in both cases (Fig. 2a, 2b) is represented as being viewed by the agent, not as being represented by the partner.

Theory of Mind

Therefore, both circuits, as depicted in Fig. 1 and Fig. 2a,b still represent egocentric systems. We will now proceed allowing the agent to be able to simulate the behaviour and the internal view, including the sensory experience, of the partner, a property that has been characterized as ToM. To this end, we will present a simple simulation of the Sally-Anne task mentioned above. To be able to represent some aspects of the memory of the partner required for this task, in our network the unit **Partner** is given a connection to memory elements representing the position of the candy as viewed by the partner (Fig. 3, dashed line). Now imagine that subject Anne is either equipped with a network as depicted in Fig. 2 or in Fig. 3. Application of a system shown in Fig. 2 means that the agent (Anne) has only one representation of the candy's position, the one seen last. Therefore only this, correct, position can be activated and the partner is imagined to grasp the correct position as observed in children younger than about four years. The child is not taking into account the position the partner assumes. In contrast in a system as presented in Fig. 3a, there is a difference in thinking of oneself grasping the candy or the partner doing it. When the agent imagines itself to grasp the candy, it would grasp to the correct and known position. If asked to simulate the internal state of the partner, as is required in the case of the Sally-Anne test, (Fig. 3b), the position connected to the partner Sally will be used and the agent would rightfully deduct that the partners grasp would be directed towards this position which is wrong, but this fact is not known by the partner. Therefore, the network shown in Fig. 3 allows for ToM, in contrast to the network shown in Fig. 2. The critical difference between both networks is that the network shown in Fig. 3 contains a separate representation of (a part of) the partner's memory. Ishida et al. (2010) describe mirror neurons that are able to represent this property.

Feeling pain

To illustrate another, more difficult case, let us come back to a push movement being directed to a partner. This case is more complex because roles can be interchanged in this scenario as the partner could also push the agent. To simulate this situation, the Ego network has correspondingly to be equipped with an element containing its spatial position, called "pos.Ego" in Fig. 4 (to simplify the figure, elements "grasp" and pos.candy are omitted in this and the later figures). In the following, two possible situations are considered, (1) the agent pushing the partner ("Ego push Partner") and (2) the partner pushing the agent ("Partner push Ego"). In these situations the agent may act as an actor (corresponding to a grammatical subject in an active phrase) or as a patient (corresponding to a grammatical object in an active phrase). Therefore, instead of having one unit for each individual as in the networks explained above, we introduce now two units to represent each individual, the agent and the partner. The corresponding subject units and object units are arranged under the column "subject" and "object" (Fig. 4a,b) and are connected via mutual inhibition.

To represent a verbally given situation like "Ego push Partner" in the network, some way is required to define roles. Here we assume that the item given first in time functions as

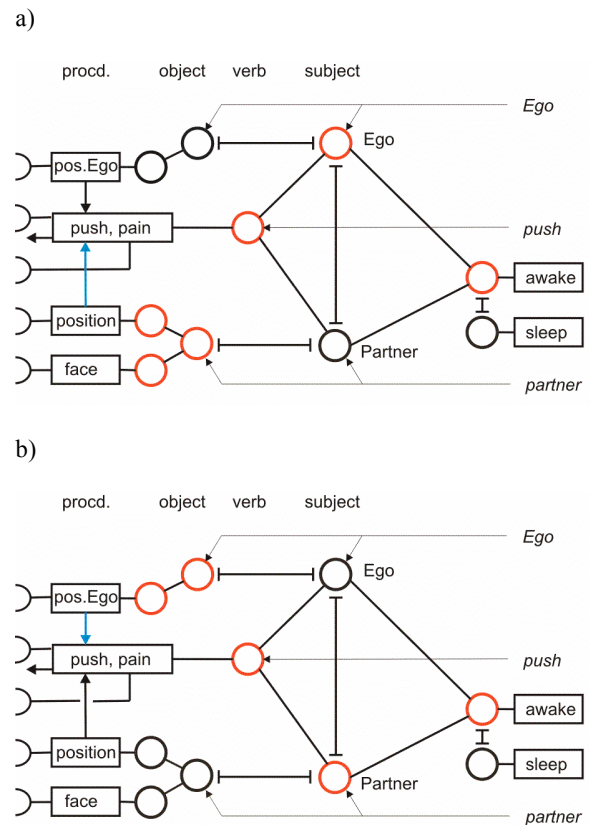


Fig. 4. A network allowing for ToM, being able to represent an egocentric view (a, situation "Ego push Partner") and the view as seen by the partner (b, situation "Partner push Ego"). Units for individuals (agent, partner) can be represented by an 'object unit' or a 'subject unit', as indicated in the top line. Sensorimotor, or procedural, networks can be found under the heading 'procd.', action units under 'verb'. For further explanation see Fig. 1 and text.

subject, the second as verb, and the third as object. The network shown in Fig. 4a,b maps the temporal order into the neuronal structure. Beginning with situation (1) input *Ego* is given first and is immediately followed by *push*. This leads to an activation of the unit *Push* and the *Ego*-subject unit (Fig. 4a, red) and an inhibition of both the *Ego*-object unit and the *Partner*-subject unit. *Ego*-subject unit is activated rather than the *Ego*-object unit because only the former is supported by activation of the unit *Push*. Later, both partner units will be activated via input "*partner*". As the *Partner*-subject unit is already inhibited, the *Partner*-object unit will win, in turn activating its position unit (Fig. 4a, red). Thus, all units required to represent situation (1)—the agent performs a push directed to the partner position—are active. In this way, this network can represent the egocentric view as was already possible for the networks shown in Figs. 1 and 2.¹

¹ If the situation is not given by verbal input, but for example by visual observation, the roles of the different items actor, action and patient may

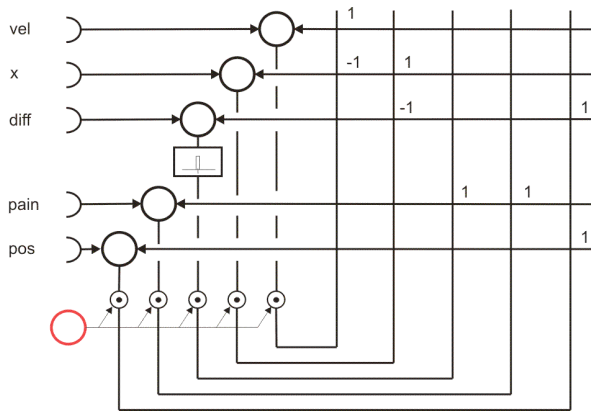


Fig. 5. A recurrent network using five IC units that shows in more detail the sensorimotor element termed “push, pain” in Fig. 5. The uppermost three units represent a simple (one-dimensional) form of the push controller (vel: velocity of the end-effector, x: position of the end-effector, also used as motor output, diff: spatial difference between actual position and goal position, represented by unit “pos”). The recurrent network “pain”, consisting of one unit, when activated long enough represents the neuronal substrate for feeling pain. The unit diff possess a nonlinear activation function that allows to activate the pain network when the activation of the unit diff has approached a value of about zero. The activation of the complete network is controlled by a motivation unit (red circle).

The same network can however correspondingly represent situation (2) “Partner push Ego”. To this end, the partner units, now representing the actor, are first activated together with Push, whereas in a later step unit Ego is activated. In a corresponding way, at the end Partner-subject unit, unit Push, as well as Ego-object unit and Ego-position unit remain active (Fig. 4b).

If the agent is confronted with the latter situation “Partner push Ego” for the first time, it may suffer from a painful feeling, which will then be associated with being pushed. The network whose activation is accompanied by the subjective experience of pain (Fig. 4, box ‘push,pain’), is integrated into the push network in the following way. The pain network is activated when the controlled position of the tip of the arm reaches the goal position, the pain being associated with the goal position.

To illustrate how the networks push and pain and the input from the position network are connected, in Fig. 5 a minimal version of this subnetwork is depicted in more detail. The network altogether consists of five IC units plus one motivation unit. The push network contains three units, one representing position of the end-effector of the arm characterized by one dimension, x, the (constant) velocity of the end-effector, vel, and a unit diff representing the difference between the actual position x and the target position pos. Unit diff has a nonlinear activation function

be internally represented by different salience values provided by neuronal systems able to detect these different roles.

providing an output of 1 in a small interval around an activation value of zero, and providing a zero output otherwise. In all three cases, one unit suffices to represent the corresponding values as we focus on a one-dimensional example.² Furthermore, there is an RNN, consisting of one unit that when activated represents a painful state (pain). Unit pain is activated as soon as the end-effector meets the target position (diff = 0). We will not deal with the question how these weights are learned.

If—after this network has been installed and the situation (1) “Ego push Partner” is activated (either as active behaviour or only as imagined, i.e. simulated, behaviour)—the position of the partner will be associated with the feeling of pain (arrow highlighted in blue in Fig. 4a). In this way, our agent can simulate and thereby experience the experience of the partner without confusion between the two individuals. This means that the agent shows the ability being endowed with empathy (following the definition of Decety and Jackson, 2004: “Empathy accounts for the [...] subjective experience of similarity between the feelings expressed by self and others without loosing sight of whose feelings belong to whom”).

Coming back again to the second situation (Fig. 4b), “Partner push Ego”, the agent can simulate the view of the partner being an actor. Now the position of the agent is provided to the push network (in Fig. 4b depicted by a blue arrow). Therefore the network of the agent can simulate that the agent himself is receiving a push and experiencing a painful feeling. Thus, the simulated partner can now be experienced as to experience the pain.

Taken together, the agent equipped with a network as shown in Figs. 3, 4 can experience an egocentric view as was already possible for the networks shown in Fig. 1 or 2 (see Figs. 3a and 4a). In addition, the agent is able to ‘put himself into the shoes of the partner’ in two ways: the agent can try to understand the view of its partner onto objects (Fig. 3) or onto itself (Fig. 3b and 4b), i.e. “seeing himself with the eyes of the conspecific” (ToM), and can experience the experience of the partner (Fig. 4a) by simulating the feeling of the partner. The simulation of the partner is of course based on the innate and learned structures underlying his own ability to feel.

Mutualism

A further evolutionary as well as developmental step that, according to Tomasello (2009) is unique for humans, is described by the term mutualism. Mutualism concerns the property of an agent to cooperate with another agent in such a way that both individuals perform—possibly different—actions by which a common goal should be reached and where both individuals will profit. A simple case is to carry a heavy load (e.g. to move a table around obstacles). A formally related task has to be solved by a hexapod walker where the legs are considered to be driven by independent controllers, but the legs being mechanically coupled via the body and the ground. For this problem two different solutions have been proposed. One solution possibly realized by insects exploits the mechanical coupling of the legs applying an extremely

² Note that we reduce these networks to a minimum size in order to better explain the essential aspects. Of course, each network could be expanded to consist of a large number of units without touching the basic statements made here.

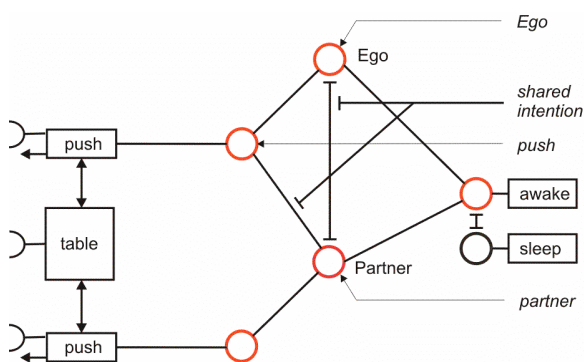


Fig. 6. A network allowing for the control of mutualistic behaviour. If input “shared intention” is activated, the (excitatory and inhibitory) connections between the subnetworks representing the agent (Ego) and the partner are interrupted. Therefore, both subnetworks can be used simultaneously to simulate actions that pursue a common goal. For further explanations see text and Fig. 1.

decentralized control structure (Schmitz et al., 2008). As an in our context more interesting alternative, Cruse and Schilling (2010) and Schilling (2011) proposed the application of an internal model that allows to simulate the legs plus their mechanical coupling through the world. Using this model each leg controller provides commands to its leg in such a way that each individual leg supports the common goal, namely moving the body forward. Applying this example to our problem of considering two independent agents able to behave mutualistically, the controller of each agent should correspondingly possess a model not only of itself, but also of the partner and the relevant environmental conditions. Together, these three elements form a ‘supermodel’. In analogy to Tomasello’s terminology, this model might also be called a “we-model”. Application of this supermodel can correspondingly be used for *probehandeln*, i.e. imagined behaviour, in order to reach a common goal. Indeed, Tomasello argued that the ability to have a we-mode is a prerequisite for developing a common goal.

What are the requirements for such a we-model to be implemented? First, the ability has to be given that actions of both the agent and the partner can be simulated independently and simultaneously. This means that it does not suffice to have only one body model that can be used to either simulate the Ego or the partner as was the case for the ‘shared-circuits’ networks shown above (Figs. 2, 3, 4). Rather, both motivation units, Ego and Partner, require access to separate behavioural elements (e.g. push) and a body model each. In Fig. 6, as in Figs. 2, 3, 4, the body model is not shown explicitly, but is graphically embedded in the push network. Both body models have to be connected via a model simulating (part of) the world to represent the actual situation, in our example the table to be carried. Furthermore, to activate the we-model, the mutual inhibition between both motivation units Ego and Partner has to be suppressed (Fig. 6). A suppression is also necessary for the connection between the motivation unit Partner and the push model which in the networks shown in Figs. 3, 4 is necessary because the latter is shared between the

Ego and the Partner network. Therefore, the we-model is activated by an input termed shared intention in Fig. 6. Tomasello has already considered shared intention a crucial property for a system showing mutualism. If this mode has been adopted, the we-model can be used to search for a solution to a given problem, for example moving the table. This search, of course, takes into account actual sensory information, e.g. position of the table relative to both individuals, movement of the other individual and possibly verbal information.

Discussion

In our earlier work, we proposed a network that is able to control behaviour (walking, climbing, navigation) using a behaviour-based architecture and that has been expanded to show a fundamental cognitive ability, namely to be able to plan ahead. Here we propose several expansions of this network, *reaCog*. As these expansions follow the basic structure of *reaCog*, they can easily be implemented in the *reaCog* architecture. Using a typical section of *reaCog*, as an example, we start with an egocentric system (Fig. 1) that contains a body model, but is not able of mirroring. In the first step, we introduce a new connection that allows the egocentric system to apply a mirror system, i.e. to interpret behaviours observed when being performed by other individuals (Fig. 2). However, application of shared circuits alone does not appear sufficient to allow for the representation of how the world is represented by others, i.e., to allow the network shown in Fig. 2 to solve the Sally-Anne task. The latter is however possible for the networks developed in the next step (Figs. 3 and 4), which in addition contain a representation of parts of the partner’s memory. The latter concerns the position of an object, the candy in the example shown in Fig. 3 or the position of the partner (Fig. 4). In the latter example, (Figs. 4, 5), we explain in more detail how this system might implement subjective feelings which could be attributed to other individuals. Both networks are able to apply ToM. The architecture shown in Fig. 4 is still based on the application of shared circuits as the push/pain network can be connected to either the unit Ego or the unit Partner. Separation into subject units and object units is required to represent the different roles the agents have to play in this paradigm. In contrast to the egocentric systems (Figs. 1, 2), the systems depicted in Figs. 3 and 4 may be called allocentric.

Fig. 6 shows what additional connections may be required to allow for mutualism. Here two body models can be activated simultaneously and the connections allowing for sharing circuits are inhibited. Application of such a we-model is suited to allow for what often has been called the third person’s view. The step from a network as shown in Fig. 4 to that presented in Fig. 6 appears to correspond to an idea proposed by Keyesers and Gazzola (2011) who draw a distinction between application of shared circuits, used for mirroring to understand the partner at a lower, intuitive, non-cognitive level, and another system involving different brain areas when subjects are asked to reflect on others. According to Keyesers and Gazzola, both mechanisms are activated according to the abstraction level of the actual task. Such a two-body model appears also to be helpful to explain a number of experimental results reviewed by Sebanz et al.

(2006) and Vesper et al. (2010) which show that subjects require shared representations of tasks including the simulation of the expected behaviour of confederates.

It might be tempting to speculate that the existence of these two body models might form the basis of some illusory own-body perceptions where, due to specific neuronal deficits, subjects can experience two body representations and self-identification refers either to the physical body (Autoscopy), to the illusory body (Out-of-Body experiences) or to both either simultaneously or in alternation (Heautoscopy) as described by Blanke and Metzinger (2009). In our system such illusions may result if accidentally both body models are connected to the unit Ego, a connection not depicted in Fig. 6.

Acknowledgements

This work has been supported by the Center of Excellence 'Cognitive Interaction Technology' (EXC 277), by the EC-IST EMICAB project # FP7 – 270182 and by a DAAD postdoctoral fellowship.

References

- Bläsing, B. (2006). Crossing large gaps: A simulation study of stick insect behaviour. *Adaptive Behavior*, 14(3):265–285.
- Blanke, O. and Metzinger, T. (2009). Full-body illusions and minimal phenomenal selfhood. *Trends in Cognitive Sciences*, 13:7–13.
- Chalmers, D. J. (1996). *The conscious mind*. Oxford University Press, New York.
- Cooper, R.P. and Shallice, T. (2006). Hierarchical schemas and goals in the control of sequential behavior. *Psychol. Rev.* 113:887–916.
- Cruse, H. (2003). The evolution of cognition — a hypothesis. *Cognitive Science*, 27(1):135–155.
- Cruse, H. (1999). Feeling our body — the basis of cognition? *Evolution and Cognition*, 5:162–173.
- Cruse, H. (2010). The talking stick: A cognitive system in a nutshell. In Giuliani, L., editor, *Jahrbuch Wissenschaftskolleg zu Berlin*, pages 52–61. Wissenschaftskolleg, Berlin.
- Cruse, H., and Schilling, M. (2010). Getting cognitive. In Bläsing, B., Puttke, M. and Schack, T., editors, *The Neurocognition of Dance*, pages 53–74. Psychology Press, London.
- Cruse, H. and Wehner, R. (2011). No need for a cognitive map: Decentralized memory for insect navigation. *PLoS. Comp Biol.*
- Decety, J., Jackson, P.L. (2004). The functional architecture of human empathy. *Behav. Cogn. Neurosci Rev.* 3: 71-100.
- Dürr, V., Cruse, H. and Schmitz, J. (2004). Behaviour-based modelling of hexapod locomotion: Linking biology and technical application. *Arthropod Struct. Develop.* 33(3):237–250.
- Feldman, J. and Narayanan, S. (2004). Embodied meaning in a neural theory of language. *Brain and Language* 89 (2):385–392.
- Freud, S. (1911). Formulierung über die zwei Prinzipien des psychischen Geschehens. In *Gesammelte Werke*, Bd. VIII, pages 229–238.
- Gallese, V. & Lakoff, G. (2005). The brain's concepts: the role of the sensory-motor system in conceptual knowledge. *Cognitive Neuropsychology* 22(3–4):455–479.
- Ishida, H., Nakajima, K., Inase, M., Murata, A. (2010). Shared mapping of own and others' bodies in visuotactile bimodal area of monkey parietal cortex. *Journal of Cognitive Neuroscience* 22:83-96.
- Jeannerod, M. (2006). *Motor Cognition — What Action tells the Self*. Oxford: University Press.
- Jeannerod, M. (1999). To act or not to act: Perspectives on the representation of actions. *Quarterly Journal of Experimental Psychology*, 52A:1–29.
- Keyesers, C. and Gazzola, V. (2011). Integrating simulation and theory of mind: from self to social cognition. *Trends in Cognitive Sciences*
- Kilner, J.M., Paulignan, Y. and Blakemore, S.J. (2003). An Interference Effect of Observed Biological Movement on Action. *Current Biology* 13:522–525.
- Kühn, S., and Cruse, H. (2007). Modelling memory functions with recurrent neural networks consisting of input compensation units: II. Dynamic situations. *Biological Cybernetics*, 96(5):471–486.
- Makarov, V.A., Song, Y., Velarde, M.G., Hübner, D. and Cruse, H. (2008). Elements for a general memory structure: properties of recurrent neural networks used to form situation models. *Biological Cybernetics* 98(5):371–395.
- McFarland, D., Bösner, T. (1993) *Intelligent behavior in animals and robots*. MIT Press, Cambridge, MA
- Menzel, R., Brembs, B. and Giurfa, M. (2007). Cognition in Invertebrates. In Kaas, J.H., editor, *Evolution of Nervous Systems, Vol. II: Evolution of Nervous Systems in Invertebrates*, chapter 1.26, pages 403–422. Academic Press, Oxford.
- Narayanan, S. (1997). Talking the talk is like walking the walk: A computational model of verbal aspect. In COGSCI-97, pages 548–553. Stanford, CA.
- Rizzolatti, G. and Luppino, G. (2001). The cortical motor system. *Neuron* 31: 889-901.
- Schilling, M. (2011). Universally manipulable body models — dual quaternion representations in layered and dynamic MMCs. *Autonomous Robots* 30(4):399–425.
- Schilling, M. and Cruse, H. (submitted). Cognition as recruitment of reactive systems.
- Schilling, M. and Spranger, M. (2010), "Embodied posture verbs: Emergence of a vocabulary for describing postures", Conference on Conceptual structure, discourse and language (CSDL) & Embodied and situated language processing (ESLP) 2010, San Diego.
- Schmitz, J., Schneider, A., Schilling, M. and Cruse, H. (2008). No need for a body model: Positive velocity feedback for the control of an 18-DOF robot walker. *Applied Bionics and Biomechanics*, 5(3):135–147.
- Sebanz, N., Bekkering, H., Knoblich, G. (2006). Joint action: Bodies and minds moving together. *Trends in Cognitive Sciences* 10:70–76.
- Spranger, M., Höfer, S. and Hild, M. (2009). Biologically inspired posture recognition and posture change detection for humanoid robots. In *Proceedings of ROBIO'09: IEEE International Conference on Robotics and Biomimetics*, pages 562–567.
- Steels, L. (1995). A self-organizing spatial vocabulary. *Artificial Life*, 2(3):319–332.
- Steels, L. (2003). Intelligence with representation. *Philosophical Transactions: Mathematical, Physical and Engineering Sciences* 361 (1811):2381–2395.
- Steels, L. and Spranger, M. (2008). The robot in the mirror. *Connection Science*, 20(4):337–358.
- Tomasello, M. (2009). *Why we cooperate*. Boston Review Book, MIT Press, MA.
- Vesper, C., Butterfill, S., Knoblich, G., & Sebanz, N. (2010). A minimal architecture for joint action. *Neural Networks*, 23, 998-1003.
- Wehner, R. (2008) The desert ant's navigational toolkit: procedural rather than positional knowledge. *Navigation* 55(2):101–114.

This document was last revised on June 6, 2011.