

Neural dynamic concepts for intentional systems

Jan Tekülve (jan.tekuelve@ini.rub.de)

Gregor Schöner (gregor.schoener@ini.rub.de)
Institut für Neuroinformatik, Ruhr-Universität Bochum
44780 Bochum, Germany

Abstract

How may intentionality, the capacity of mental states to be about the world, emerge from neural processes? We propose a set of theoretical concepts that enable a simulated agent to have intentional states as it perceives, acts, memorizes, plans, and builds beliefs about a simulated environment. The concepts are framed within Dynamic Field Theory (Schöner et al., 2015), a mathematical language for neural processes models at the level of networks of neural populations. Inspired by Searle’s analysis of the two directions of fit of intentional states (Searle, 1980), we recognize that process models of intentional states must detect the match of the world to the mind (for “action” intentions) or the match of the mind to the world (for “perceptual” intentions). Neural representations of Searle’s condition of satisfaction implement these detection decisions through dynamic instabilities that are instrumental in enabling autonomous switches among intentional states.

Keywords: Dynamical systems modeling; Mathematical modeling; Neural networks; Intelligent agents; Cognitive Architectures

Introduction

How are neural processes organized to create coherent, complex cognitive function? For instance, how are sequences of actions and processes of active perception generated to orient actions at objects to achieve a desired outcome? How may the nervous system switch between actions and mental states that are driven by current sensory information and actions or mental states that are driven by memory and knowledge?

Philosophers of mind have framed related questions in terms of the notion of *intentionality*: How may an organism with its nervous system generate intentional states that are about objects in the world? How may an organism act to change the world according to its intentional states? The logical structure of this problem has been analyzed in depth by John Searle (Searle, 1980). He postulates that intentional states come in two directions of fit (DoF), the *world-to-mind* direction of fit, in which an intentional state’s content represents a desired state of the world, capturing the intuitive “action” flavor of intention. The *mind-to-world* DoF comprises states in which the state’s content matches circumstances in the world, a “perceptual” flavor of intention. Each intentional state can be described through its *condition of satisfaction* (CoS), which determines whether the fit between mind and world is achieved. Searle has conjugated these two forms of intentionality through three layers of psychological modes: *intention-in-action* (IiA) and *perception* are intentional states

directly linked to the motor or sensory systems. *Prior intention* and *memory* are intentional states with a more indirect form of linkages, in which additional steps are needed to act out or bring about the intentional state. *Beliefs* and *desires* are more abstract forms of intentionality, typically thought to take propositional forms, with an inherent generalization beyond the immediately accessible perceptual or motor experience.

We come to these questions from the theoretical framework of Dynamic Field Theory (DFT) (Schöner et al., 2015), a mathematical language for neural processes models at the level of networks of neural populations. Here, we take inspiration from Searle’s concepts to address the neural processes required to autonomously switch between intentional states in these six psychological modes. A key idea has been that there must be neural processes that explicitly represent a CoS and whose activation controls the transitions from one intentional state to another (Sandamirskaya & Schöner, 2010). Specifically, for world-to-mind intentional states, activation of the neural representation of the CoS signals the successful achievement of an intentional state that leads to its deactivation and opens the system to switch to a subsequent intentional state. In mind-to-world intentional states, it is the representation of the content of the intention itself that forms the CoS, which is activated when a detection decision is made and remains activated as long as the intentional state persists.

In this paper we develop this idea into a systematic account of how intentional states can be organized to autonomously generate goal- and object-oriented behavior. We simulate a rudimentary toy scenario, in which an agent explores its simple environment containing colored objects and buckets of paint. The agent may move towards objects and direct an effector to them, either taking up paint (for a bucket) or painting the object (for the colored objects). The agent detects objects, may attentionally select objects, may build scene memories, generate sequences of actions to paint particular objects with a particular paint, and learn and exploit beliefs about which paint applied to which surface generates which outcome. Simple desires (to seek particular outcomes of painting acts) drive the agents goal-oriented and exploratory behaviors. The scenario is chosen such that the amount of time each action or mental operation takes varies, and that during that time the agent is exposed to other perceptions or sensory states that could distract from its current intention. The inherent stability of its intentional states and the capacity to

release these states from stability under the right conditions is thus probed in this scenario.

Dynamic Field Theory

Dynamic Field Theory (DFT) (Schöner et al., 2015) is a theoretical framework for understanding perception, motor behavior, and cognition based on neural principles. The activity in neural populations is modeled by activation fields, $u(x, t)$, spanned across the metric dimensions, x , to which the population is tuned. The neural dynamics of the activation fields,

$$\tau \dot{u}(x, t) = -u(x, t) + h + s(x, t) + \int \omega(x - x') \sigma(u(x', t)) dx$$

describes the time-continuous evolution of neural activation on the time scale τ . Activation $u(x)$ below the sigmoidal threshold σ relaxes to the stable solution $h + s(x)$, defined by the field's resting level h and its localized inputs $s(x)$. Field sites, where activation strength surpasses the threshold level, will engage in lateral interaction defined by the field's kernel $\omega(x - x')$, which is locally excitatory and inhibitory over longer distances $x - x'$. This leads to the formation of self-stabilized peaks of supra-threshold activation, which are the unit of representation in DFT (see figure 1).

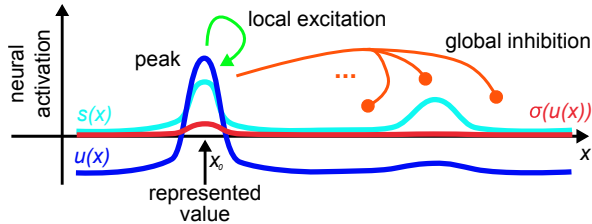


Figure 1: A dynamic neural field spanned across the metric-dimension x representing value x_0 through a supra-threshold activation peak.

Depending on the individual strength of excitatory and inhibitory interaction, fields may allow the formation of multiple peaks (self-stabilized), single peaks (selective) or they may sustain peaks once localized input is removed (self-sustained). Multi-dimensional fields may represent conjunctions of feature dimensions, for example, the conjunction of color and space. Zero-dimensional fields are dynamic neural nodes that represent categorical states.

Two fields u_{src} and u_{tar} may be coupled by adding a field's output $\sigma(u_{\text{src}})$ to the other field's rate of change \dot{u}_{tar} , weighted with a homogeneous connection kernel $\omega_{\text{src, tar}}$. Such projections may preserve the dimensionality of the fields, or may expand or contract the field dimensionality (Zibner & Faubel, 2015). *Dimensionality expansions* may take the form of ridges (or tubes, or slices), in which input along one or several of the receiving field's dimension is constant. *Dimensionality contractions* typically entail integrating along the contracted dimension. Dynamic neural nodes that project homogeneously onto a field by expansion are called *boost nodes*. They may alter the dynamic regime in the target field

and induce the formation or vanishing of peaks. Within field architectures such boost nodes may effectively modulate the flow of activation by enabling or disabling particular branches of an architecture to create units of representation. *Concept nodes* project a specific pattern on a higher dimensional field to elicit a peak representing the concept, e.g. a blue-concept node activates neurons tuned to blue hue in a field spanned across the color dimension.

The transition from a stable sub-threshold solution to a new supra-threshold activation pattern marks a discrete event in the presence of time-continuous input variations and is labelled *detection instability*. In the context of intentional states the detection instability is utilized to determine a state's condition of satisfaction, the discrete point in time where a successful match between world and mind representations is achieved.

In the world-to-mind DoF a *matching field* (CoS field) receives sub-threshold input from an *intention-field* representing the desired world-state and sub-threshold input from a *perception-field* representing the current world-state. Due to the resting level h in relation to the strengths of both field inputs, a supra-threshold peak will only form in the matching-field, if both input patterns overlap sufficiently, thus signaling the states CoS through a detection instability. Representation of a world-to-mind CoS is thus independent from the planned timing of the underlying action and signals its termination on a perceptual basis. The formation of a CoS may thus be used to terminate the action and activate the next action in a planned sequence (Richter, Sandamirskaya, & Schöner, 2012).

In the mind-to-world DoF the CoS is determined through the formation of a peak in a field that is connected to sensor or memory substrates. The detection instability may be the result of salient input alone or of the combination of sensor/memory input and top-down attention input from within the neural architecture. Representations of a mind-to-world CoS are made available to the rest of the architecture and may be used in further cognitive processing, e.g. determining a world-to-mind CoS.

Transforming Searles logical analysis of intentional states into a process account has led us to a number of new insights. One is a difference in the time structure of world-to-mind vs. mind-to-world intentional states. World-to-mind intentional states are active before the corresponding state of the world has been achieved and are deactivated once the CoS detects a match between the expected and the sensed state of the world. Mind-to-world intentional states, in contrast, often persist beyond the detection of a match, which is an essential characteristic of memories and beliefs. But what if memories or beliefs (and even percepts) are false? Then they must be deactivated. This is controlled by a condition of dissatisfaction (CoD), which detects a mismatch between current sensory or internal information and an intentional state. Upon activation, a CoD inhibits that intentional state. The CoD responds to evidence against the intentional state, not to the mere absence

of evidence supporting the intentional state.

Model/Scenario

We illustrate how intentional states can be organized to generate autonomous goal- and object-oriented behavior in a minimal scenario requiring Searle’s six major psychological modes. The scenario contains a simulated agent engaged in an artificial painting task controlled by a dynamic field architecture connected to the robots sensorimotor surface (see figure 2 for a sketch).

Mind-to-World States

Intentional states of the mind-to-world DoF are the prerequisite to engage in meaningful actions in a given environment as any action at least aims to achieve a perceivable outcome.

Perception The virtual environment contains cuboids of different height and color, which are arranged in an array along a single dimension facing the robotic agent. The agent’s visual perception fields are therefore spanned across horizontal retinal space and the two feature dimensions height and color. A selective spatial attention mechanism causes peaks to form in the same spatial location in the *space/color* and

space/height perception fields, representing a perception of the particular height and color features at that particular location (see Grieben et al. (2018) for details on the attentional selection). To detect successful interaction with the world, the agent perceives changes in the environment through a two-layer transient detector that forms peaks in response to sudden changes in visual input (see Berger et al. (2012) for details).

To monitor its own actions the agent requires self-perception of the task-dependent “body parts”, which includes an estimate of the agent’s position in the world. A simulated sensor provides input to a one-dimensional *current position* field, as the agent’s movement is restricted to driving in parallel to the cuboid array. Arm movement is restricted to two Cartesian dimensions, lateral and forward translation, which leads to a two-dimensional representation of the current end effector position in the *proprioception* field. The painting device is located at the robot’s end effector and can either be filled with color or not. This categorical perceptual state is represented through a neural node that is activated if the device is filled.

Attention directed towards particular self-perceptions is modeled through a homogeneous resting level boost, which causes the sub-threshold sensor information to form a peak

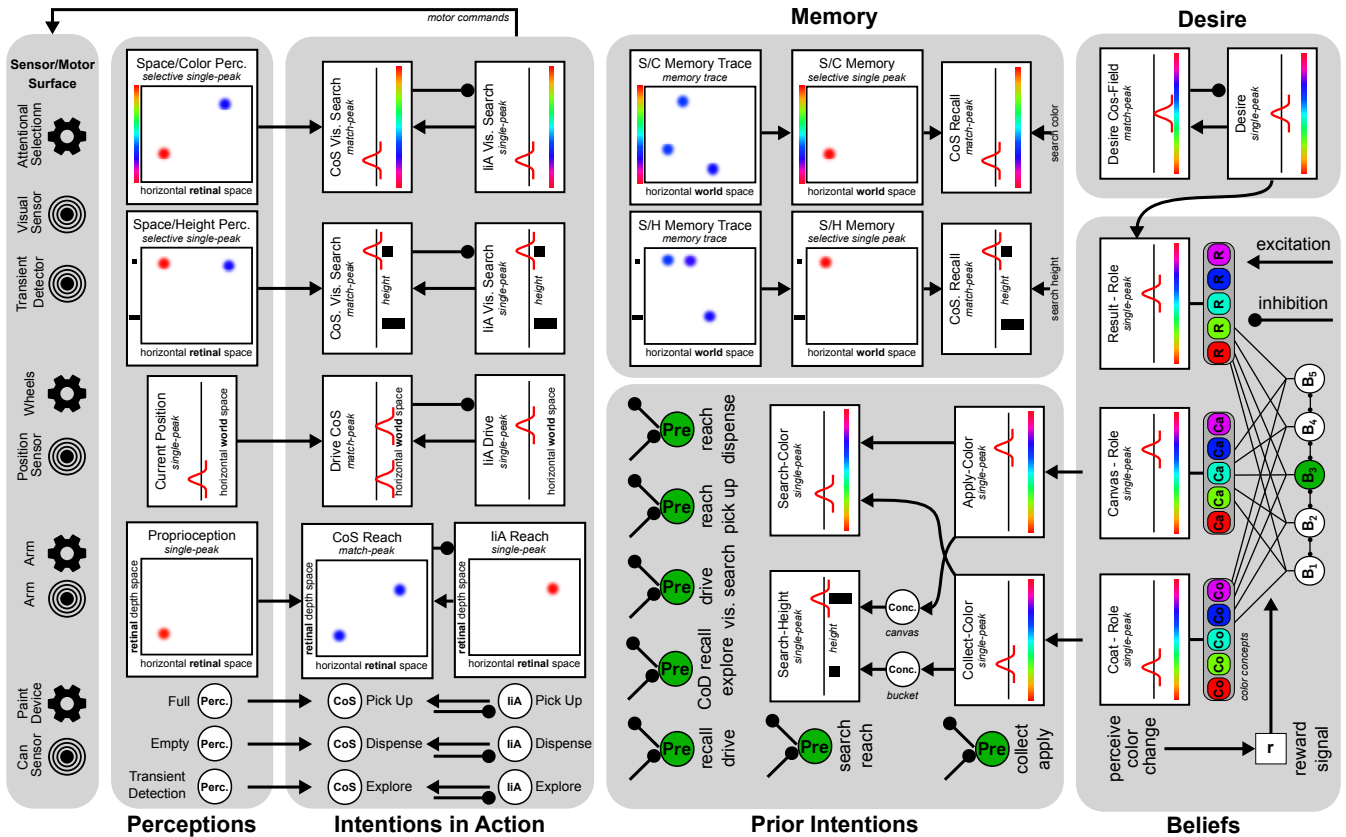


Figure 2: Schematic overview of the dynamic fields and nodes representing the agent’s intentional states grouped according to their psychological modes. For clarity’s sake only the most relevant connections are shown and parts of the architecture relevant to autonomous learning and exploration are hidden. Prior intentions are depicted as precondition nodes with labels describing the inhibiting CoS followed by the inhibited IIA.

in the respective perception field. Neural interaction in perception fields is strong enough to prevent the destabilization of perceptions through noise, but retains its input coupling such that a continuous change in input induces a drift in peak position.

Memory To allow the agent to engage in more sophisticated actions that are not purely based on current perceptions, the agent stores past perceptions of cuboids in memory. Each visual perception of the agent leaves a slowly decaying two-dimensional memory-trace spanned across world-space and feature, modeling a memory process that is subject to interference (Erlhagen & Schöner, 2002). The trace is forwarded as sub-threshold activation to a space/feature *memory* field analog to the visual perception fields. Memory states represented as peaks in the memory field may emerge through either spatial or feature cues overlapping with the memory trace substrate.

Self-sustained fields retaining task-relevant information, such as the recently collected color, represent working memory, which is functionally closer to the mode of perception than memory, as self-sustained peaks resemble lasting perception representations and do not need an additional detection mechanism to form.

Belief Meaningful interaction with the world also relies on general knowledge or beliefs about the world represented in propositional form. In the toy scenario, beliefs are about relations between the three color concepts: the color of a canvas, the color of the paint, and the color that results from coating the canvas with the paint. Each painting action contributes to the formation of a belief about that relation. The relation is represented through a neural node with reciprocal connections to three color concept nodes, each linked to a different *color role* field. An activated belief state is represented through a supra-threshold belief node that leads to the formation of three peaks, each in one self-sustaining color role field, which provide working memory representations to guide the painting process. The color concept nodes ensure a degree of generalization, as different shades of hue activate the same concept node, while the activation of the concept node activates the mean hue value of the particular color.

A belief is activated when color nodes in either of the three roles become active, to which the belief has learned synaptic connections. For instance, a belief linking the red point on a blue canvas to a yellow result may become activated, if the result color node yellow is activated by a corresponding desire. Inhibitory coupling among belief nodes ensures that only a single belief may be activated at any time. The belief with most matching color role input will typically win the competition and can then be used to guide action. If an active color role does not match the learned projections of any belief, no belief is activated.

The learning of new beliefs is organized by a neural dynamic architecture inspired by Adaptive Resonance Theory (Carpenter & Grossberg, 2016) illustrated in Figure 3. It as-

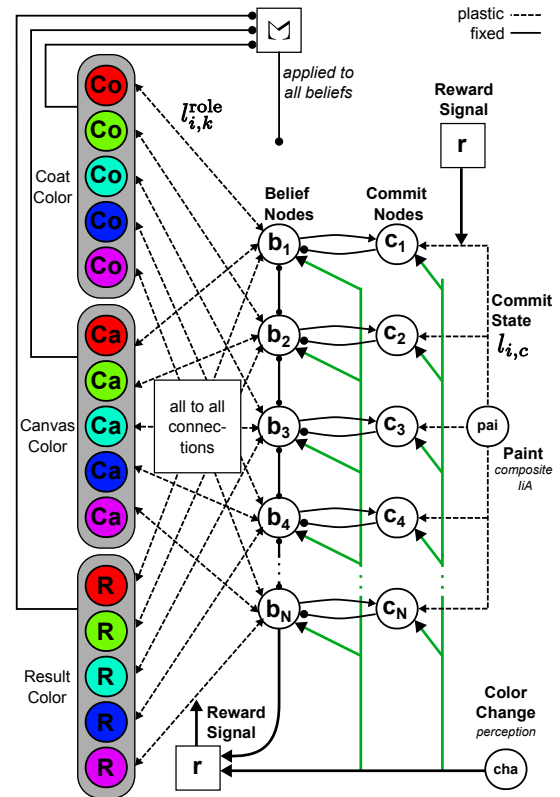


Figure 3: Detailed sketch of the belief learning architecture shown in Figure 2.

sociates color concepts in the three roles coating, canvas and result color with a single belief node. Such learning steps occur whenever the transient detector registers a change of color during a painting action. This happens under two possible conditions. In one case, a belief has previously been activated that predicts the expected color change. If that prediction is confirmed, a Hebbian learning mechanism consolidates the connectivity. If that prediction is not confirmed, the CoD is activated, and the belief is inhibited. This leaves the system without any activated belief. That second case, no activated belief, may also arise because there was no matching belief to begin with. In this case, a belief node is recruited for learning the new association between coating, canvas, and resulting color. This happens through a homogeneous boost of all belief nodes. Only a previously uncommitted belief node has a chance to become activated, because each belief node is inhibited by a dedicated “commit node” that represents that this belief node is committed to a particular belief it has learned.

The actual learning processes is modulated by a transient reward signal, $r(t)$, that is generated in the presence of an active belief node and a detected color change in the scene. The reward modulated Hebbian learning rule adapts the connections, $i_{i,k}^{\text{role}}$, between belief nodes, b_i , and color-role concept nodes, u_k^{role} (where k is color and $\text{role} \in \{\text{coat}, \text{canvas}, \text{result}\}$):

$$i_{i,k}^{\text{role}} = \eta r(t) \sigma(b_i) \sigma(u_k^{\text{role}}).$$

The learning rate, η , is chosen such that a new belief is learned within a single transient epoch of reward in a form of one-shot learning. For a more detailed analysis of the mechanisms of autonomous learning see (Tekülve & Schöner, 2019).

World-to-Mind States

Intentional states of the world-to-mind DoF are instrumental in bringing about a desired state of affairs in the world, which includes the agent’s own body. All world-to-mind states share the representation through the pair of intention and CoS (or match) field. Outgoing connections from the intention field specify the actions driven by a supra-threshold activation peak, while the outgoing inhibitory connection from the CoS field terminates the action once the desired state is detected.

Intention in Action The painting scenario provides several elementary actions that may take a variable amount of time and thus require a representation of a CoS to verify their successful execution. *Reaching* to a particular location in the visual array is realized through a neural field architecture for generating arm movements (see Zibner et al. (2015)). Its duration depends on the relative distance between the agent and the target location. The target location is defined through spatial input from the visual perception fields, which classify reaching as an object oriented action.

Moving to a particular position in the world is motivated through memory instead of perception. The *drive* IiA field thus receives its spatial input from peaks formed in the space/feature memory fields. In absence of a particular target location the agent may also move to either direction until a previously unattended cuboid is perceived. The *explore* IiA realizes this behavior and its CoS is represented through a binary neural node receiving excitatory input from the visual transient detector. The actions *explore*, *pick up* and *dispense* represent a family of IiAs, where the desired world state is categorical and represented through the activation of a neural node.

Another family of IiAs is represented by the actions *visual search*, *recall* and *activate belief*, which treat the current state of the neural system as part of the world and try to induce particular states of the mind-to-world DoF. Visual search guides the attentional system to achieve a perceptual state matching an intended feature cue, while recall tries to achieve a memory state matching an intended feature cue and activate belief intends to activate a belief node that matches certain color-roles.

Prior Intention Most goal directed actions comprise a sequence of actions such as the painting task in this scenario which requires: Searching for a “color bucket” (high cuboid), collecting color from it, searching for a “canvas” (small cuboid) and applying the collected color on it. Those actions themselves may be described as sequences of more elementary actions, e.g. searching comprises the sequence of recalling a cuboid’s position, driving to the position and visually

searching for the cuboid, while collecting and applying comprise reaching followed by picking up or dispensing color.

Such a sequence of actions (or composite IiA) is realized through an intention-field that simultaneously activates all IiAs involved in the sequence and an inhibiting precondition node for each IiA. The combination of activating and specifying the input of an IiA, while simultaneously inhibiting it, represents a prior intention. The prior intention turns into an IiA once the precondition node is destabilized by the CoS of a preceding IiA which releases the IiA-field from inhibition (Richter et al. (2012)). These CoS fields may sustain activation in a working memory representation of the current stage within a sequence.

The CoS of a composite IiA is activated through a subset of CoS representations of comprising IiAs determining the successful completion of the composite IiA’s goal. This will inhibit the composites IiAs intention field and subsequently destabilize all working memory representations of the comprising IiAs, thus allowing the same sequence of actions to be activated again, which is required in the scenario as searching for a cuboid is part of both collecting and applying color.

Prior intentions may also represent alternative action plans that may occur when a precondition node is destabilized by a CoD, for example, due to failing to recall a specific cuboid.

Desire The agent’s desire to observe the change of a cube’s color into a desired color is the drive for all actions it executes. The desire specifies the agent’s prior intentions of collecting and applying color through the activation of a belief that matches the desired result color. The *desire CoS* is activated through a match between a changing color detected by the visual transient detector and the desired color, which leads to a subsequent inhibition of the desire returning the field architecture to its initial state.

Results

Figure 4 shows activation snapshots of selected fields displaying the formation of CoS peaks during a successful painting sequence. In snapshot t_1 the desire to paint a cube yellow feeds into the result-role field (left column), which triggers a detection instability in belief node B_4 leading to a complete belief representation through the emergence of peaks in the canvas and coat role fields (right column). The coat color leads to an activation of the collect IiA to retrieve blue color and a prior intention to apply the color to a purple canvas, which is represented through a sub-threshold peak in the apply IiA field.

At t_2 the IiA collect activates the “bucket” concept, a high cuboid, which is forwarded as a recall cue to the space/color and space/height memory fields respectively. The collect color also forms the prior-intention to visually search for blue color (left column). The color/height cue leads to the emergence of a single memory peak at the location of the blue/high cuboid, which is read out across space and leads to the formation of a peak in the IiA drive (right column).

The left column of snapshot t_3 shows the IiA drive-field

ory). The agent follows these plans by driving towards objects and interacting with them (intention-in-action), if they visually match specified feature combinations (perception). Beliefs about the three different color roles are learned autonomously during each painting sequence.

Similar goals are pursued by Schrodt and colleagues (Schrodt & Butz, 2016; Schrodt et al., 2017), who learn production rules within a cognitive architecture. That work is framed within a probabilistic approach, which is partially embedded in neural networks. Our methods to achieve autonomous sequencing overlap with techniques developed in (Kazerounian & Grossberg, 2014). Globally speaking, we pursue similar aims as the research program of cognitive architectures (Anderson, 1996). Our emphasis is to be pervasively consistent with neural principles, generating the sequence of processing steps autonomously from neural dynamics alone. Although the functions fulfilled by portions of the neural dynamics can be described using concepts of information processing, the system is simply a set of integro-differential equations that generate time courses of activation. These integro-differential equations capture the time-continuous evolution of activation in populations of cortical and subcortical neurons (Erlhagen, Bastian, Jancke, Riehle, & Schöner, 1999). It remains a challenge to provide direct neural support for a complex model like ours (see (Wijeakumar, Ambrose, Spencer, & Curtu, 2017) for an outline of how that may happen). Empirical support for a model like ours may also be sought in the form of behavioral signatures of the neural dynamics, an approach that has been successful for past DFT models. The highly integrative nature of the model makes this difficult, but perhaps not impossible.

Future modeling tasks include scaling the demonstrated principles to more complex task-environments, elaborating the simplistic account for desires, and addressing how believed propositions may be both true and false.

In conclusion, we have explored the requirements on neural processes that arise when embodied cognitive systems are endowed with intentional states of the two directions of fit and the six psychological modes that provide a foundation for intentionality.

Acknowledgments

Support by the Deutsche Forschungsgemeinschaft (SPP Active Self, SCH 336/12-1) and by the Studienstiftung des Deutschen Volkes is gratefully acknowledged.

References

Anderson, J. R. (1996). Act: A simple theory of complex cognition. *American Psychologist*, 51(4), 355.

Berger, M., Faubel, C., Norman, J., Hock, H., & Schöner, G. (2012). *The counter-change model of motion perception: An account based on dynamic field theory* (Vol. 7552 LNCS).

Carpenter, G. A., & Grossberg, S. (2016). *Adaptive resonance theory*. Springer.

Erlhagen, W., Bastian, A., Jancke, D., Riehle, A., & Schöner, G. (1999). The distribution of neuronal population activation (DPA) as a tool to study interaction and integration in cortical representations. *Journal of Neuroscience Methods*, 94(1), 53–66.

Erlhagen, W., & Schöner, G. (2002). Dynamic field theory of movement preparation. *Psychological Review*, 109(3), 545–572.

Grieben, R., Tekülve, J., Zibner, S. K. U., Schneegans, S., & Schöner, G. (2018, July). Sequences of discrete attentional shifts emerge from a neural dynamic architecture for conjunctive visual search that operates in continuous time. In J. Z. Chuck Kalish Martina Rau & T. Rogers (Eds.), *Cogsci 2018* (pp. 427–432).

Kazerounian, S., & Grossberg, S. (2014). Real-time learning of predictive recognition categories that chunk sequences of items stored in working memory. *Frontiers in Psychology*, 5, 1–28.

Richter, M., Sandamirskaya, Y., & Schöner, G. (2012). A robotic architecture for action selection and behavioral organization inspired by human cognition. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on* (pp. 2457–2464).

Sandamirskaya, Y., & Schöner, G. (2010). An embodied account of serial order: How instabilities drive sequence generation. *Neural Networks*, 23(10), 1164–1179.

Schöner, G., Spencer, J. P., & DFT Research Group, T. (2015). *Dynamic thinking: A primer on dynamic field theory*. Oxford University Press.

Schrodt, F., & Butz, M. V. (2016). Just imagine! learning to emulate and infer actions with a stochastic generative architecture. *Frontiers in Robotics and AI*, 3, 5.

Schrodt, F., Kneissler, J., Ehrenfeld, S., & Butz, M. V. (2017). Mario becomes cognitive. *Topics in cognitive science*, 9(2), 343–373.

Searle, J. R. (1980). The intentionality of intention and action*. *Cognitive Science*, 4(1), 47–70.

Tekülve, J., & Schöner, G. (2019). Autonomously learning beliefs is facilitated by a neural dynamic network driving an intentional agent. In *Ieee conference on development and learning and epigenetic robotics (icdl-epirob 2019)*.

Wijeakumar, S., Ambrose, J. P., Spencer, J. P., & Curtu, R. (2017). Model-based functional neuroimaging using dynamic neural fields: An integrative cognitive neuroscience approach. *Journal of Mathematical Psychology*, 76, 212–235.

Zibner, S. K. U., & Faubel, C. (2015). Dynamic scene representations and autonomous robotics. In *Dynamic thinking: A primer on dynamic field theory* (pp. 227–246). Oxford University Press.

Zibner, S. K. U., Tekülve, J., & Schöner, G. (2015). The neural dynamics of goal-directed arm movements: a developmental perspective. In *Ieee conference on development and learning and epigenetic robotics (icdl-epirob 2015)* (pp. 154–161).