

Autonomy of 'Self' at Criticality: The Perspective from Synthetic Neuro-Robotics

Jun Tani

RIKEN Brain Science Institute

2-1 Hirosawa, Wako-shi, Saitama 351-0198, Japan

Tel: +81-48-467-6467, FAX: +81-48-467-7248

E-mail: tani@brain.riken.jp

Abstract

The current paper investigates the phenomenological aspects of “selves” in relation to autonomous agents. Through a review of a series of neuro-robotics experiments conducted by the author’s group, we elucidate three different aspects of “selves”, namely, minimal selves, social selves and self-referential selves. Upon integrative discussions of these “selves”, it is suggested that genuine constructs of “authentic” selves may appear with criticality, which is self-organized in the iterative interplay between regression of past experience and lookahead prediction of future outcomes. It is concluded that genuine autonomy of agents likely originates from genuine autonomy of authentic selves.

1 Introduction

The boom in the work on autonomous agents started two decades ago with the publication of the edited book, “Designing Autonomous Agents” (Maes, 1991). The concepts of autonomous agents presented in this book differed drastically from those considered in conventional artificial intelligence research, which had emphasized abstraction with explicit symbolic representation and deliberative “thinking” processes for planning and inference. These new studies focused on the bottom-up pathway from sensorimotor interactions rather than the top-down internal “thinking” process. One key idea in this new paradigm was “emergent functionality” (Maes, 1991). The notion was that, even with simple rules governing sensorimotor interactions, quite complex and unpredictable behaviors could emerge in the coupled dynamics between the agents and their environments.

It is certainly interesting to watch some behavior-based robots built on such ideas (Brooks, 1991) interact with their environments. One such example is a mobile robot exploring light sources while skillfully avoiding obstacles. Another is a robot participating in simple interactions with humans such as stopping in front of us and saying “hello” or following us (Horswill, 1993). Such interactions sometimes stimulate us to think about agency for these robots. However, after a while, we may begin to feel that the robots with reflex behaviors are simply like steel balls in pinball machines, repeatedly bouncing against the pins until they finally disappear down the holes; while we may recognize some complexity on the surface level of these behaviors, they are fundamentally different from those generally expected from humans. Therefore, we are likely at some point to become bored with interacting with such robots. Those robots that have passed the Turing test (Turing, 1950) turn out to be just machines having stochastic state transition tables. But what is wrong with these robots? Although they have neither complex skills for action nor complex concepts for conversation, such complexity issues may not be the main problem.

The current paper conjectures that the problem originates from a fundamental lack of phenomenological constructs in those robotic agents. In particular, what is missing might be the “subjectivity” that should direct their intentionality to project their own particular images on the outer objective world. Such subjectivity should be developed gradually through interactions with direct sensorimotor experiences with the world. Development of own views or particular internal models would enable robots to anticipate and to interpret the outcomes of their actions.

My argument for the need for subjectivity is related to the issue of identity in

defining agency, as expounded by Barandiaran, Paolo and Rohde (Barandiaran, Paolo, & Rohde, 2009). They consider that a necessary condition for the appearance of agency is the presence of a system capable of defining its own identity as an individual and thus distinguishing itself from its surroundings. Of particular interest here is their view that the identity or boundary of individuals should be self-defined through environmental interactions. And the same should also hold true for the concept of subjectivity in my discussion. It is important to note that anticipation when seen from the top-down view of agents may not fit with real-world reality in many situations. When environmental interactions proceed exactly as expected, behaviors can be generated smoothly with automaticity. However, anticipation can sometimes be wrong, and the conflict that arises in such case can make generating successive acts difficult. Furthermore, when the state of the interactions shifts spontaneously between these opposite poles of the automaticity with perfect matching and the struggle with the conflicts, the boundary between the subjective mind and the objective world would fluctuate. Here, I argue that the essential characteristics of this phenomenon would be better understood by revisiting discussions of “selves” in traditional phenomenology since phenomenologists have already investigated the dynamic characteristics of the “autonomy” of selves in their language. Thus, I attempt here to situate or anchor the problems of agency in the literature of phenomenology.

In the next section, we examine in more detail phenomenological accounts of selves by discussing some of the relevant literature. We then review a series of our neuro-robotics experiments, whose results might correspond to some aspects of phenomenological selves. Lastly, we attempt to postulate a definition of the organizational principle of autonomy of selves in humans as well as in robotic agents by developing triangular discussions from the perspectives of nonlinear dynamics, neuroscience and phenomenology, following the so-called neuro-phenomenological approach proposed by Varela (Varela, 1996).

2 Phenomenological selves

We humans are self-conscious beings. However, the state of self-consciousness is elusive as we experience it every moment of our lives. We sometimes strongly feel the existence of our 'selves' as isolated in the world, but the feeling does not last for long. Or, we can be apart from such feeling of 'selves' for hours when concentrating on certain tasks, but the feeling can return suddenly. William James, the pioneering American psychologist and philosopher, wrote that, when we take a general view of the wonderful stream of

our consciousness, what strikes us first is the different pace of its parts. Like a bird's life, it seems to be an alternation of flights and perchings (James, 1890). Our group believe that autonomy of cognitive systems might originate from just such dynamics of spontaneous shifting between substantive phase of resting and transitive phase of flying in a stream of consciousness.

Some phenomenologists consider that self in the substantive phase is like a concrete object. Strawson (Strawson, 1997) suggests the image of a string of pearls as a metaphor of self, claiming that each self should be considered as a distinct particle-like existence, an individual thing or object, like a pearl, yet discontinuous as a function of time. Then, from where does this sort of discreteness or particle-like characteristics of self come from, and how can self be consciously aware? Phenomenology has found clues in the interactions between 'subject' and 'object' ¹. Let us take Heidegger's (Heidegger, 1962) well-known example of the hammer. For a carpenter, when everything is going smoothly, the carpenter himself and the hammer function as a unity. But, when something goes wrong with the carpenter's hammering or with the hammer, then the independent existences of the subject (the carpenter) and the object (the hammer) are noticed. Here, the carpenter becomes conscious of himself as problematic, just as he becomes conscious of the world, due to things not going as expected. The substantive phase in the stream of consciousness may correspond to the unity where the subjective inner reality in terms of expectation creates a perfect match with the objective reality, whereas the transitive phase may correspond to the breakdown of unity when self becomes consciously aware of the gap. In other words, momentary self exists as a discrete aspect of unity but its existence is noticed only at the moment of its breakdown ².

Here, we might ask how he or she can regard all the fragmented discrete selves as actually belonging to her or his coherent self. Gallagher (Gallagher, 2000) regards momentary subjective experience of self as the 'minimal self'; however, with reference to Hume (Hume, 1975), he claims that the minimal self can be developed to the narrative self constituted with a past and a future in the various stories that we tell about ourselves.

This construction of the narrative self as like stories from subjective momentary experiences of self might be deeply related to Husserl's thoughts on time perception

¹Although it is apparently nontrivial to define what are subject and object in phenomenology, these terms are introduced here for explanatory purposes. It will be described later that such distinction exists merely in external observer domains.

²There is a recent alternative discussion that minimal selves can be constituted without actions, but only with passive multi-sensory body images (Blanke & Metzinger, 2009).

(Husserl, 1964). An essential phenomenological question on time perception voiced by Husserl (Husserl, 1964) concerns how objective time can be constituted out of the subjective flow of temporality. Although temporality is experienced as a part of flow at a deep phenomenological level, it is experienced as temporal objects and events at a shallow level. In recalling past experiences, a temporal image of the past can be reactivated as a linear sequence of discrete events, rather than as a replay of the original continuous flow of our impression.

Varela (Varela, 1999) pointed out an apparently paradoxical nature in the human perception of temporality, in that a conscious event is perceived as a unity but it is still part of a flow. He proposed the use of nonlinear dynamics systems theory as a formal descriptive tool for the phenomenon. By using the phenomenon of spontaneous flipping in multi-stable visual perception, such as in the Necker-cube illusion, he explains that the dynamic properties of chaos, which is characterized by spontaneous shifts between static and rapid transition modes, could explain the paradox of continuous, yet segmented time perception.

The same sort of phenomenon is also recognized to occur in action generation. It has been reported that an action slip, or failure of an action pattern referred to as a micro-slip (Reed & Schoenherr, 1992), appears in our everyday actions. For example, while making a cup of coffee, one may mistakenly grasp the coffee cup rather than the spoon. An interesting observation here is that the action slip does not occur randomly during the course of the entire action but at some meaningful point. More specifically, the entire action of, for example, making a cup of coffee seems to consist of meaningful chunks such as grasping the coffee cup or grasping the spoon where the slip can take place not inside the chunks themselves but at segmentation or branching points between the chunks. The chunks here might correspond to behavior primitives or motor schemata in Arbib's motor schemata theory (Arbib, 1981), in which primitives or schemata are considered as a set of reusable action units. Therefore, it is considered that the stream of ongoing motor behaviors ought to have complementary properties of constancy and flow (Varela, 1999) where constancy might be afforded by convergent dynamics trapped by accumulated motor schemata and flow might be afforded by divergent dynamics of freely combining them.

In regards to the construction of the narrative self, from bundles of experiences of pure sensorimotor flow, the flow in its original form cannot be manipulable unless it is somehow segmented into a set of identifiable objects. Then, the ultimate question is how subjective experiences of continuous sensorimotor flow can be transformed into manipulable objects by which selves can be consciously described. In other words, how

can subjective experiences of selves be reflected by themselves? This question actually addresses the issue of self-referential selves (Varela, 1999; Butz, 2008).

On the basis of this question, by closely examining the dynamic phenomena observed in our synthetic robotics experiments, the current paper attempts to postulate that neuronal processes of self-organization are likely the underlying key mechanism constituting self-referential selves. The following review of these experiments will illustrate that possible constructs for self-referential selves can emerge in the internal neuronal dynamics through self-organizing compositional mechanisms of assembling and de-assembling sensorimotor schemata of repeated experiences. And, most importantly, it will be suggested that self-referential selves could be constituted only in *critical* conditions of sustaining conflictive interactions between the top-down 'subjective' mind and the bottom-up sensorimotor reality.

Before closing this section, I will consider the social aspect of selves. A curious question here is how we can recognize the selves of others. This question is considered to be related to agency detection problems, which were addressed by Trevarthen's double TV-monitor experiments (Trevarthen, 1993) as well as Auvray's perceptual crossing paradigm (Auvray, Lenaya, & Stewart, 2009). A particularly interesting finding in the double TV-monitor experiments is that mother and infant can engage in coordinated utterances and affective expressions when 'live' facial interactions are allowed on TV monitors but the infants cannot so engage when the mothers appear in video-recorded interactions on the monitor (Trevarthen, 1993). Auvray's perceptual crossing paradigm showed the same analogy in adult studies. Some model studies (Iizuka & Paolo, 2007; Martius, Nolfi, & Herrmann, 2008) simulating Auvray's findings have postulated that social interaction may lead to agency detection of each other when mutual anticipation is formed dynamically. In another line of research, Wolpert et al. (Wolpert, Doya, & Kawato, 2003) postulated that anticipatory competency for others by means of mirror neurons (Rizzolatti, Fadiga, Galless, & Fogassi, 1996) can lead to agency detection in the theory-of-mind context. However, our essential question regarding these considerations is how we can distinguish one-self and others if one can perfectly anticipate about others through a unity which is formed among them by coherent coupling. Although mutual interactions should be inevitable for agency detection, perfect coherence may not be necessary prerequisite. Genuine perception of agencies or selves of others should originate primarily from their unpredictable parts and incoherences, as will be postulated later in this paper.

The next section will start to describe our synthetic robotics approach that would address the issue of autonomy of selves and agency.

3 The dynamical systems approach

My colleagues and I have utilized the dynamical systems approach (Schoner & Kelso, 1988; Beer, 1995; Tani, 1996) in building our cognitive robots. An advantage of this approach is that it can deal with continuous sensorimotor space directly, whereas other formula, such as symbol systems or probabilistic modeling, require arbitrary partitioning of the original continuous sensorimotor space into a set of discrete states. A feature of our models that is distinct from most other agent models utilizing the dynamical systems approach is that our models have top-down pathways based on neural activation patterns which interact with the bottom-up pathways based on sensorimotor patterns. These two pathways share the same metric space and thus their interactions can be dense (Tani, 1996). While there have been some studies (Holland & Goodman, 2003; Ziemke, Jirenghed, & Hesslow, 2005) on consciousness and selves that have installed top-down pathways of anticipation, these studies did not necessarily pay close attention to their dynamic interactions with the bottom-up pathways. In the light of this, our ultimate expectation is that close investigation of the characteristic nature of such dynamic interactions between the two pathways might shed light on the continuing phenomenological discussions about the relationship between the subjective mind and the objective world.

Let us now review three different robotics experiments that we conducted on the basis of this notion. The first experiment may account for appearances of 'minimal selves' in a simple robot navigation experiment, the second for appearances of 'social selves' in an imitation learning experiment between robots and human subjects, and the third for appearances of 'self-referential selves' in a more complex skill learning experiment. In this review, descriptions will be highlighted with observed phenomena related to problems of selves, but not for details of the adopted models or experimental setups which the reader may find in (Tani, 1998) for Experiment-1, (Ito & Tani, 2004) for Experiment-2 and (Tani, 2003) for Experiment-3.

4 Experiment-1: 'Minimal Selves'

A mobile robot with a vision camera mounted on a rotating head shown in Figure 1 (a) learns to predict landmark sequences it encounters dynamically, while the robot navigates in a closed workspace. After successful learning, the robot should be able to visually attend to landmarks of colored objects and corners in appropriate timing before encountering them while traveling around the workspace by following the wall and

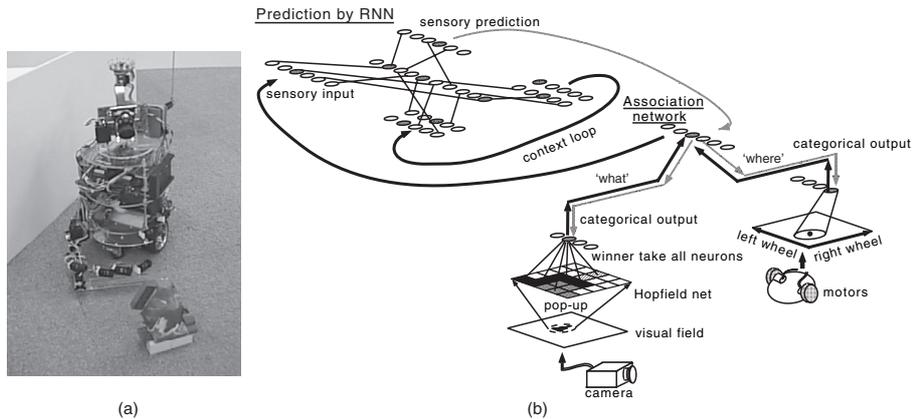


Figure 1: (a) The vision-based mobile robot looks at a colored landmark object. (b) Neural network architecture employed in the robot.

visually attending to the boundary between the wall and the floor.

4.1 Model and experiment setup

The robot is controlled by the neural network architecture shown in Figure 1 (b). The entire network consists of prediction by a recurrent neural network (RNN) (Jordan, 1986) and perception, which is divided into the 'what pathway' and the 'where pathway'. In the what pathway, visual patterns of the attended landmarks of colored objects are processed in a Hopfield network (Hopfield & Tank, 1985). When perception of a visual pattern converges into one of the learned fixed point attractors, the pattern is recognized and its categorical output is generated by a Kohonen network (Kohonen, 2001). The learning is conducted for both the Hopfield network and the Kohonen network whenever a visual stimulus is encountered. In the 'where pathway', accumulated encoder readings of left and right wheels from the last encountered landmark to the current one and the direction of detected landmarks in the front view are processed by a Kohonen network and its categorical outputs are generated. Together with both pathways, 'what' categories of visual landmark objects and 'where' categories of relative travel distance from the last landmark to the current one as well as its direction targeted by the camera head are sent for prediction in a bottom-up manner.

For prediction, a Jordan-type RNN (Jordan, 1986) learns to predict the perceptual categories of 'what' and 'where' for landmarks to be encountered next in a top-down manner. In this model, the bottom-up and the top-down pathways do not merely provide inputs and outputs to the system. Rather, they exist for their mutual interactions.

The system is prepared for expected perceptual categories in the top-down pathway before actually encountering the landmarks. This expectation ensures readiness for the next arriving pattern in the Hopfield network as well as readiness to direct the camera head toward the landmark with correct timing and direction. Actual perception is established by dynamic interactions between the two pathways. This means that if the top-down prediction of the visual pattern does not match the currently encountered one, the perception would result in an illusion of a combined pattern of the two. Moreover, mismatch in the 'where' perceptual category can result in failure to find any of the expected landmarks. Such perceptual outcomes are fed into the RNN and the next prediction is made based on this. Note that the RNN should not be considered as just an input/output mapping system. The implementation of so-called context units recurrency (see Figure 1 (b)) makes the RNN an autonomous dynamical system where the inputs act as perturbations in the ongoing dynamics. At this point, the RNN is regarded as an implementation of Maturana and Varela's idea (Maturana & Varela, 1980) that a neuronal network is a 'closed circuitry' without inputs and outputs. Note also that the RNN can generate 'mental rehearsing' of learned sequential images by feeding the current prediction outputs back to the next inputs, a process referred to as closed-loop operation.

For the purpose of achieving adequate interactive balance between the top-down and the bottom-up pathways, a particular mechanism for internal parameter control is implemented. The mechanism exerts more top-down pressure on the both perceptual categories of 'what' and 'where' as the error between the predicted perception and its actual outcome is lessened. A shorter time period is also allocated for reading the perceptual outcomes in the Hopfield network in this case. On the other hand, when the error between the predicted perception and its actual outcome is larger, less top-down pressure is exerted and a longer time period is allowed for the dynamic perception in the Hopfield network. In other words, in the case of fewer errors, top-down prediction dominates the perception with quick attention to coming expected landmarks which results in quick convergence in the Hopfield network. Otherwise, the bottom-up pathway dominates the perception, taking a longer period to look at landmarks while waiting for convergence in the Hopfield network.

The learning of the RNN is conducted for event sequences of landmark encountering. More specifically, experienced sequences of the perceptual category outcomes are utilized as target sequences to be learned. The incremental learning of the RNN is conducted every 15th landmark encounter by adopting a scheme of 'rehearsing' and of 'consolidation', the details of which appear in (Tani, 1998).

The experiment was conducted in a closed workspace containing five landmarks (two colored objects and three corners). The experiments were repeated three times, in each of which the robot traveled for 105 times of landmark encountering (completing approximately 21 circuits of the workspace).

4.2 Experimental results

For each run, we observed three characteristic features of the robot’s travels: prediction error, bifurcation process of the RNN dynamics due to its iterative learning, and phase plots representing attractor dynamics of the RNN at particular times in the bifurcation process. A typical example is shown in Figure 2 (a).

The prediction error was quite high in the beginning of all three trials because of the initially random connective weights. After the first learning period, the predictability was improved to a certain extent in all three trials, but the errors are not minimized completely. Prediction failures occurred intermittently during the course of the trials. We can see from the bifurcation diagram that the dynamical structure of the RNN varies from time to time. In a typical example shown in Figure 2 (a), a fixed point attractor appearing in the early period of learning iterations as a single point is plotted at each step in the bifurcation diagram mostly before the 3rd learning period. After the 3rd learning period, a quasi-periodic or weak chaotic region appears. Then, after the 4th learning period, it becomes a limit cycle of period 5 as can be seen from the 5 points plotted at each step during this period in the bifurcation diagram. In addition, its snapshot is seen in the phase plot where 5 points are plotted. After the 5th learning period, a region of strong chaos appears, as indicated by a strange attractor in the corresponding phase plot. However, the strange attractor (chaos) and the limit cycle attractor of period 5 alternate with each other. We observe that limit cycling dynamics with a periodicity of 5 appear most frequently in the courses of the all trials. The periodicity of 5 is significant since it corresponds to the 5 landmarks which the robot encounters in a circuit of the workspace. It should be noted that the observed limit cycling dynamics with a periodicity of 5 do not remain stationary; the periodicity of 5 disappears intermittently and other dynamical structures emerge.

On the basis of these results, we can surmise that there exist two distinct phases: the steady phase represented by the limit cycling dynamics with a periodicity of 5, and the unsteady phase characterized by non-periodic dynamics. It is also seen that the shifts between these two phases take place arbitrarily in the course of the time development. Moreover, it was observed that the differences appear in the physical

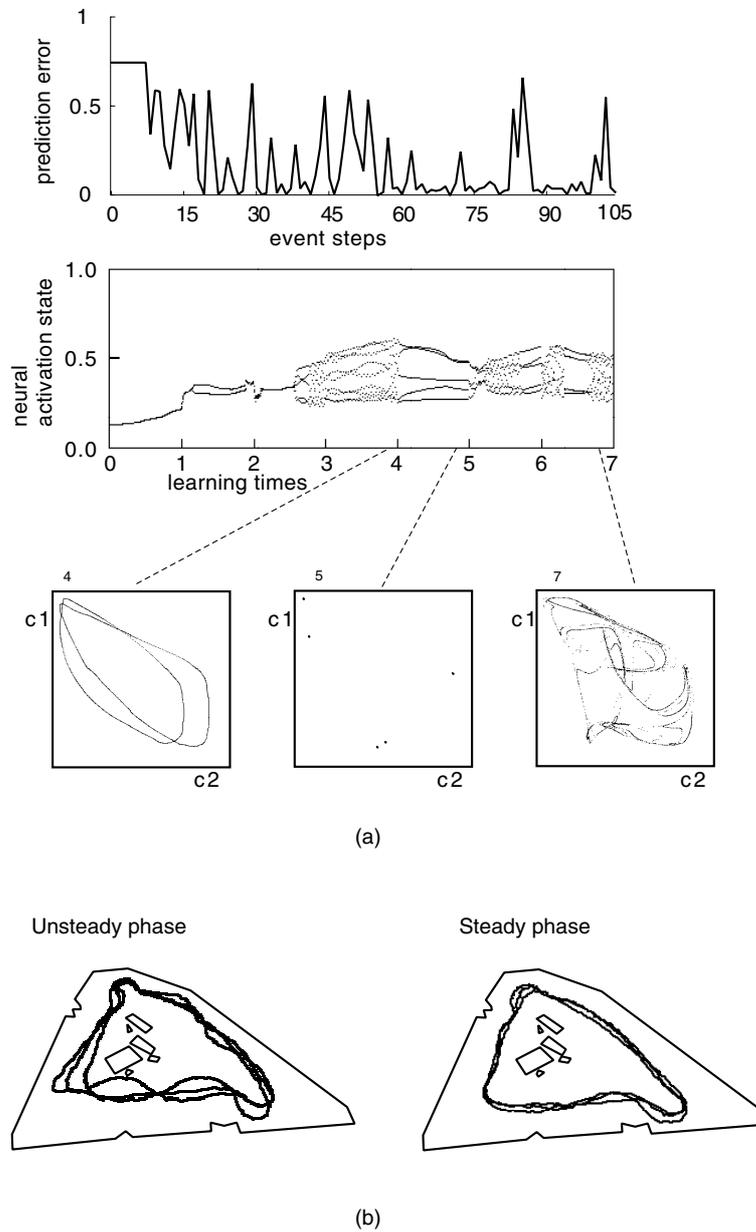


Figure 2: (a) Prediction error, bifurcation diagram of RNN dynamics, and phase plot of two context units at particular times during robot exploration learning. (b) The robot's trajectories as measured in the unsteady phase and in the steady phase.

movements of the robot as well. In order to elucidate this observation, we compared the actual robot trajectories observed in these two periods. Figure 2 (b) shows the robot trajectories measured in these two periods by a camera mounted above the workspace. The trajectory winds more in the unsteady phase than in the steady phase, particularly in the way objects or corners are approached. It is inferred that the maneuvering of the robot is more unstable in the unsteady phase because the robot spent a longer period on the visual recognition of objects due to the higher value of the prediction error. Thus, the robot faced a higher risk of misdetecting landmarks when its trajectory meandered during this period. Indeed, it misdetects corners and objects when its trajectory meandered severely during this period. In the steady phase, however, the detection sequence of landmarks became more deterministic and travel was smooth with greater prediction success. Of importance is the observation that the steady and unsteady dynamics are attributed not only to the internal cognitive processes arising in the neural network, but also to the physical movements of the robot’s body as it interacted with the external environment.

Finally, we measured the distribution of interval steps between the catastrophic error peaks ($error > 0.5$) observed in three different travels of the robot (Figure 3).

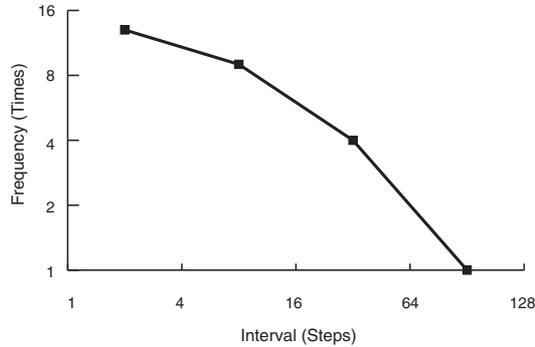


Figure 3: Distribution of interval steps between the catastrophic prediction error peaks of greater than 0.5, where the x-axis represents ranges of interval steps and the y-axis represents the frequency of appearances in the corresponding range with log scales for both axes.

The graph indicates that the distribution of the breakdown interval has a long-tail characteristic with near power-law profile. This indicates that the shift from the steady phase to the unsteady phase takes place in an unpredictable manner without dominant periodicity.

4.3 An account of 'minimal self'

An interesting observation of our experiment was that the steady phase and the unsteady phase switch with each other spontaneously even though the workspace environment was constant. Therefore, it can be said that the “inner world” of the robot is not constant at all. In the steady phase, good coherence is achieved between the internal dynamics and the environmental dynamics when subjective anticipation agrees well with observation. All the cognitive and behavioral processes proceed smoothly and automatically; no distinction can be made between the subjective mind and the objective world. In the unsteady phase, this distinction becomes rather explicit as the conflicts are generated between what the subjective mind expects and the outcome generated in the objective world. Consequently, it is in this moment of incoherence that the “self-consciousness” of the robot arises, where the system’s attention is directed to the conflicts to be resolved. On the other hand, in the steady phase, ‘self-consciousness’ is diminished substantially since there are no conflicts demanding the system’s attention. This interpretation of the experimental observations corresponds to the notion of self-consciousness in Heidegger’s (Heidegger, 1962) hammer example or Gallagher’s (Gallagher, 2000) minimal self as described in the previous section.

However, a question still remained: why could the coherence in the steady phase not last and the breakdown into incoherence take place intermittently? It seems that complex time development of the system emerges from mutual interactions between multiple local processes. It was observed that the change in the visual attention dynamics due to the change in predictability caused drifts in the robot’s maneuvering. These drifts resulted in misrecognition of the upcoming landmarks, which led to re-adaptation of the dynamic memory stored in the RNN and a consequent change in the predictability. The dynamical interactions take place between all of the processes of attention, prediction, perception, learning and behavior, and a certain *criticality* might be built up among them.

This can be explained more intuitively as follows. When the learning error decreases as the learning proceeds, more strict timing of visual attention to coming landmarks is required since only a short period of attendance to the objects is allowed proportional to the amount of the current error. In addition, the top-down image for each coming landmark pattern is shaped into a fixed one without variances. This is because the same periodic patterns are learned repeatedly and the robot trajectories tend to repeat exactly in the steady phase. If all goes completely as is expected, this strictness grows with further decreasing of the prediction error. Ultimately, at the peak of strictness,

catastrophic failures in the recognition of landmark sequences can occur even as the result of minor noise perturbation since the entire system has evolved too rigidly by building up relatively narrow and sharp top-down images.

The described phenomena might remind the reader of a theoretical study conducted on sand pile behavior by Bak and colleagues (Bak, Tang, & Wiesenfeld, 1987). In their simulation study, grains of sand are dropped onto a pile, one at a time. As the pile grows, its sides become steeper, eventually reaching a critical state. At this very moment, just one more grain would trigger an avalanche. They found that although it is unpredictable exactly when the avalanche will occur, the size of the avalanches is distributed according to a power law. This pile's natural growth to a critical state is termed 'self-organized criticality' (SOC) and it is found to be ubiquitous in various phenomena such as earthquakes, volcanic activity, the game of life, landscape formation and stock markets (Bak, 1996). Let us consider a possible analogy of the robot experiment results with the SOC phenomenon: the strictness or predictability might correspond to the steepness of the sides of the sand pile. In our robot experiment, as the strictness increases gradually with improved predictability, the system state approaches the critical state where breakdown takes place all of a sudden. In fact, our analysis has shown that the distribution of the breakdown interval has a long-tail characteristic with near power-law profile. Although we might need a larger experimental dataset to confirm SOC in the observed phenomena, our speculation is that some dynamic mechanisms for generating 'criticality' could account for the autonomous nature of 'momentary self' which William James metaphorically spoke of as an alternation of flights and perchings in a bird's life.

5 Experiment-2: 'Social Selves'

Our next experiment which concerns the imitative interactions between a humanoid robot and a human was designed to elucidate some of the social aspects of 'selves' and agency.

5.1 Model and experiment setup

In this experiment, a humanoid robot made by Sony learns to imitate multiple movement patterns demonstrated by the experimenter's hand movements (see Figure 4).

A neural network model, the so-called recurrent neural network with parametric



Figure 4: The experimenter interacting with a Sony humanoid robot.

biases (RNNPB) (Tani, Ito, & Sugita, 2004) shown in Figure 5, is utilized. The RNNPB

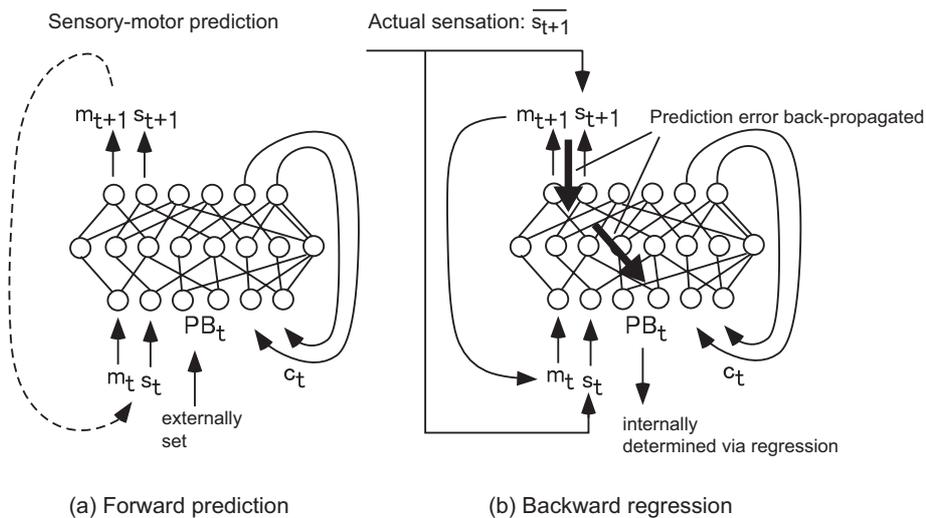


Figure 5: RNNPB model (a) forward prediction and (b) backward regression.

learns to anticipate how the human experimenter's hands change in time and also how its own arm postures should change in imitative ways through direct supervised training. The positions of the experimenter's hands are sensed by the robot visually tracking colored balls in his hands. In the interaction phase, when one of the learned movement patterns is demonstrated by the experimenter, the robot arms are expected to move by following the pattern. When the hand movement pattern is switched from one to another, the robot arm movement pattern is expected to switch correspondingly.

This sort of dynamic switching among a set of memorized spatio-temporal patterns is made possible by utilizing an online regression-prediction mechanism in RNNPB.

In terms of the RNNPB mechanism, it learns multiple sensorimotor schemes as embedded distributedly in a single RNN by utilizing additional units called parametric biases (PB). PB function as bifurcation parameters for the forward dynamics realized by the RNN. By modulating the values of the PB vectors, the forward dynamics generates diverse dynamic patterns in terms of sensorimotor sequences (s_t, m_t) by going through successive bifurcations. The essential idea here is that the PB as a bifurcation parameter works as a switcher among multiple memorized dynamic patterns. The learning in RNNPB is regarded as a process of determining an optimal synaptic weights matrix to embed all target dynamic patterns and a set of PB vectors specific to each of the target dynamic patterns. As the result of learning, mapping between the PB vector and dynamic patterns is self-organized. This can be performed by back-propagating (Rumelhart, Hinton, & Williams, 1986) the prediction error in the output units into synaptic weights as well as into the PB units. Both the synaptic weights and the PB activation values are updated in the delta error direction (see (Tani et al., 2004) for details of the implementation.) In the interaction phase, the forward prediction and the backward regression are iterated repeatedly without changing the synaptic weights. In the forward prediction shown in Figure 5 (a), the PB values, which had been self-adapted in the previous regression, as well as the current sensorimotor values (s_t, m_t) are fed into the input layer and then the next time step sensorimotor values (s_{t+1}, m_{t+1}) are predicted. In the backward regression shown in (b), the prediction error generated in the immediate past window of N steps are back-propagated into the PB units and then the PB values are updated in the direction of minimizing the error.

In the following subsection, the fundamental mechanism of the RNNPB in a master-slave type imitation setting is briefly described, and this basic experiment is then extended to include mutual interactions between human subjects and robots in which 'selves' in the social context will be discussed.

5.2 Experiments

(A) Basic experiment

After the robot was trained with 4 different cyclic patterns, we examined how it can follow shifting of the learned patterns demonstrated by the human experimenter. In repeated trials, it was observed that the robot arm movements adaptively follow each shift made by the human experimenter. Moreover, the PB vector is modulated in a

stepwise fashion at each moment of shift.

The underlying mechanism of the step-wise modulation of the PB can be explained as follows. When prediction of the human hand movements agree with their reality, there exists no error to modulate the PB values. However, when the human experimenter suddenly changes the current movement pattern to another learned one, a prediction error is generated by modulation of the PB values toward the direction of minimizing the error. This causes the robot shift its movement pattern to the correct one currently being demonstrated by the experimenter. Here, the prediction error is the drive to segment continuous sensorimotor flow into sequences of learned primitives or chunks. It could be said that 'self-consciousness' might arise at the moment when conflicts are elicited between the top-down expectation and the bottom-up reality. In other words, it is suggested that direct experience of continuous sensorimotor flow is decomposed into sequences of identifiable 'objects' by accompanying momentary consciousness of them.

(B) Mutual imitation game

The basic experiment described above involved master-slave type one-directional interaction where only the robot side adapts to the human master side, and was developed to examine mutual interaction by introducing a simple game between the robot and human. In this new experiment, after the robot learns 4 movement patterns in the same way as described previously, subjects who are unaware of what the robot learned are faced with the robot. The subjects are asked to identify as many movement patterns as possible which they and the robot can synchronize together by going through exploratory interactions. Five subjects participated in the experiment. The settings of the network and the robot were exactly the same as those in the previous interaction experiments. Each subject was allowed to explore interactions with the robot for one hour. Although most of the subjects could eventually identify all of the movement patterns, the exploration processes were not trivial for the subjects. If they merely attempted to follow the robot movement patterns, they could not converge in most instances since the PB values fluctuated when receiving unpredictable subject hand movement patterns as the inputs. If the subjects attempted to execute their desired movement patterns regardless of the robot movements, the robot could not follow them unless the movement patterns of the subjects corresponded to those learned by the robot.

An example of the interaction in the imitation game is plotted in Figure 6, where both hands positions of the subject are shown in the upper plot, their prediction by the robot is shown in the middle plot, and the 4-dimensional PB vector in the lower

plot. It can be seen that certain coherence in terms of synchronization between the

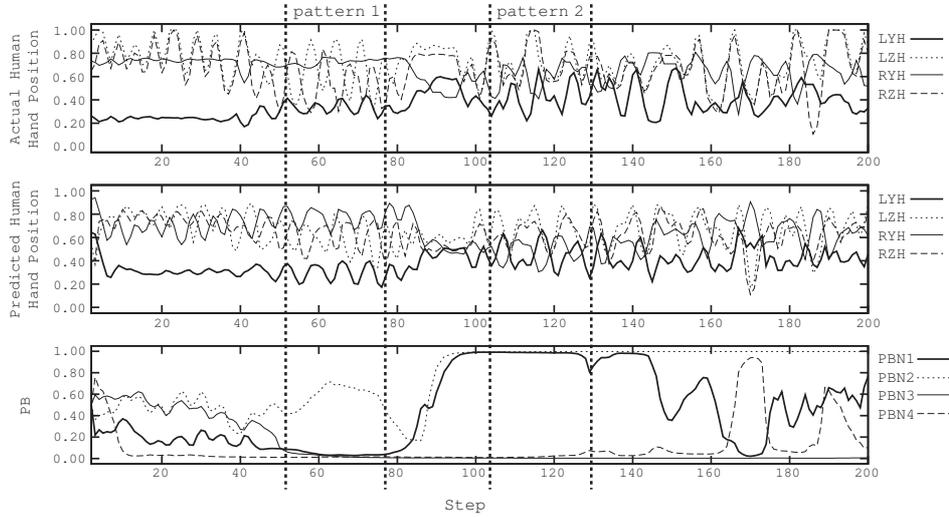


Figure 6: Actual sensation of human hand position and its prediction by RNNPB. The periods denoted by pattern 1 and pattern 2 show certain synchronization between the two.

human subject’s movements and their prediction by the robot is achieved after some exploratory phase (see patterns denoted ‘pattern 1’ and ‘pattern 2’ in the figure). It was, however, often observed that such coherence was likely to break down before coherence with another pattern was achieved again.

An interesting observation is that the master-slave relation, which was fixed between the subjects and the robot in the previous experiments, was spontaneously switched in the current experiment. In the post-experiment interviews, the subjects reported that when they felt that the robot movement patterns became close to theirs, they just kept following the robot movement patterns passively in order to stabilize the patterns. However, when they felt that their movement and the robot’s movement could not synchronize, they often initiated new movement patterns, hoping that the robot would start to follow them and become synchronized. This might account for the dynamic mechanism of turn-taking (Beebe & Lachmann, 1988; Iizuka & Ikegami, 2004).

Another interesting observation is that autonomous shifts between the coherent phase and the incoherent phase tended to occur more frequently in the middle of each session when the subject was familiarized with the robot responses to some extent. When the subjects happened to find synchronized movement patterns, they tended to

keep the attained synchronization for a moment in order to memorize the patterns. However, this coherence could break down after a while due to various uncertainties in the mutual interactions. Even small perturbations in the synchronization could confuse the subjects if they were not yet fully confident of the repertoire of the robot's movement patterns. Moreover, the subject's explorations of new movement patterns made it difficult for the robot to predict and follow their movements. It is highly speculated that such observed complexity originates from the mutually interactive mechanisms of regressing past and predicting future.

However, these complex situations were rarely seen in the early stage and the late stage of each trial session. In the early stage, subjects were unable to predict the robot movements and both the subjects and the robot proceeded with complete incoherence due to immaturity in learning. In the late stage, however, the subjects become able to predict robot movements well, where one-directional master-slave type interactions appeared. The subjects tended to make the robot patterns shift among four learned patterns at their will, as was seen in the previous experiment. It is in the middle period that subjects make good as well as bad predictions by chance after they become familiar with some parts of the patterns. Most of the subjects reported that they occasionally felt as if the robot had its own 'will' because of the spontaneity in the generated interactions in this period. This could be explained again by 'criticality' that can emerge only at a specific period with an adequate balance between predictability and unpredictability in the course of the subjects' developmental learning in the mutual imitation game. If we can predict some of the behaviors of others accurately, we may feel them to be a part of ourselves. However, if there remain some unpredictable aspects to their behavior, we may perceive agency or their own selves. Ultimately, it is argued that the theory of mind may not be about knowledge for predicting others' behaviors but rather about 'metalevel' knowledge that others' behaviors can be predictable at times and unpredictable at others.

6 Experiment-3: 'Self-Referential Selves'

This study was originally conducted for the purpose of investigating effective neuronal mechanisms for learning complex goal-directed actions by attaining compositional structures internally. More specifically, our questions concerned how the continuous sensorimotor flow experienced can be articulated into sequences of reusable behavior primitives stored in memory and how these primitives can be recombined to generate desired goal-directed actions. Although there have been some architectural proposals

utilizing different neurodynamics schemes, including the hierarchical gating network architecture (Tani & Nolfi, 1999; Haruno, Wolpert, & Kawato, 2003) and neural network model utilizing multiple time scale dynamics (Paine & Tani, 2005; Yamashita & Tani, 2008), this section reviews our prior studies (Tani, 2003) which extended the scheme of the RNNPB to incorporate hierarchical structure. It will be shown that the problem of behavior compositionality is highly related to the question of how subjective experiences can be objectified with the appearance of structures of 'self-referential selves'.

6.1 Model and experiment setup

The basic notion of extending RNNPB with hierarchy is described here, but the reader is referred to (Tani, 2003) for a detailed explanation. When the RNNPB learns to reconstruct the continuous sensorimotor sequences experienced as training targets, the PB vector values tend to change in a stepwise fashion at the moment of shift in behavior primitives in the sequences. Then, the higher level RNN, which is newly introduced, learns to predict how and when the PB vector in the lower level RNNPB changes by observing them (see Figure 7 (a).) An important point concerning implementation is that each PB unit receives additional force other than the back-propagation delta error which pressures the PB unit value to be modulated toward either extremes of 0.0 or 1.0. Then, the higher level recognizes the exact moment of segmentation as, at the least, one PB unit value flips by going across its value of 0.5. The higher level predicts the next timing of when this bit flipping in the PB vector takes place and to which bit patterns the PB vector changes. In behavior generation in the physical environment, the higher level RNN attempts to predict how the PB vector will change in time based on the learning thus far (see Figure 7 (b).) This prediction of the PB vector extends not only to the immediate future but also to the immediate past. In fact, the forward prediction sweeps for a temporal window ranging from immediate past to future iteratively by preserving temporal PB vector sequences in the window memory buffer as shown in Figure 8. Meanwhile, the PB vector in the immediate past is also modulated by the bottom-up regression if an error appears in the lower level sensory prediction. Consequently, the PB vector sequence in the immediate past window is determined by means of a balance between the pressures from the top-down forward prediction and from the bottom-up regression, and such balance is attained through iterative interactions between the two sides. It is, however, noted that experimenters should set a parameter k which determines the ratio of arbitrating the pressures from

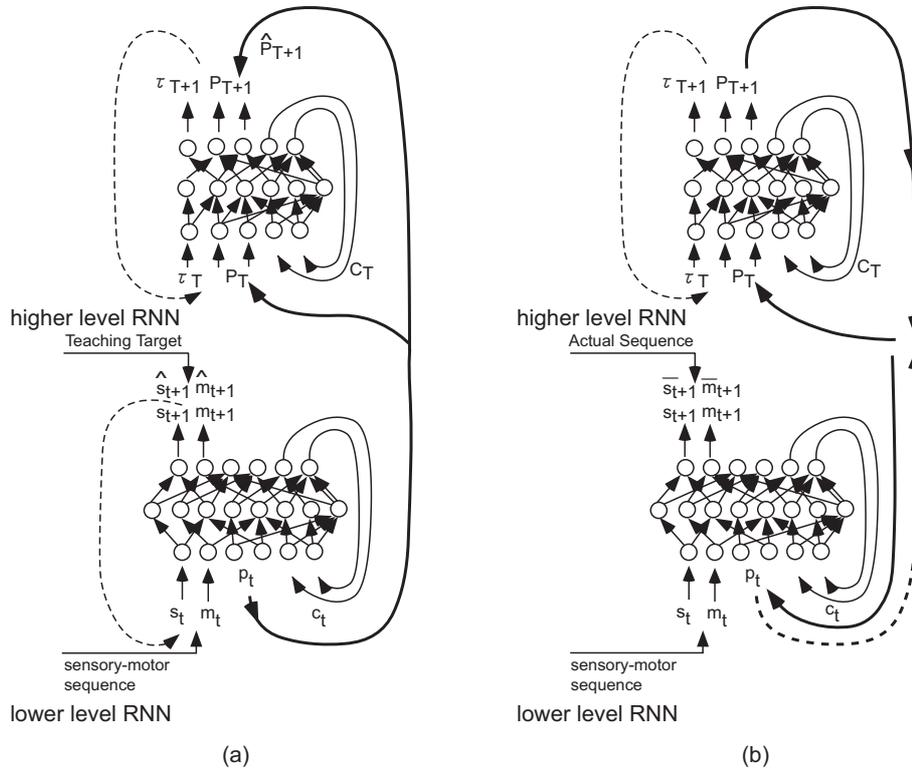


Figure 7: The RNNPB with hierarchy in which (a) the higher level RNN is trained using PB vector sequences generated in the lower level as the teaching target and (b) real-time interaction between the top-down prediction of PB from the higher level and the bottom-up regression of PB from the lower level is performed during actual behavior generation.

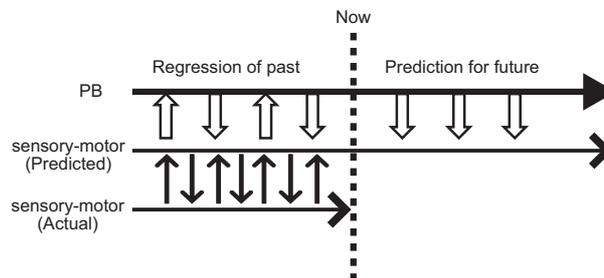


Figure 8: Prediction of future can be made by regression of sensorimotor sequences in the immediate past window.

the two sides. The actual balance during the interaction largely depends on the setting of this parameter (Tani, 2003). On the other hand, the PB vector sequence in the immediate future window is determined solely by the top-down forward prediction in the higher level. The reader may ask why this complicated mechanism is necessary. It is necessary because it is considered that future perspectives are built based on recognition of the actual behavioral outcomes in the past.

Robot experiments were then conducted for the task of simple object manipulation using a 4-degrees of freedom arm robot. In this task, the robot perceives with the robot camera the center position of the object and the arm tip position as the sensory inputs. The hierarchical RNNPB with two levels is utilized to learn multiple behavior tasks during supervised teaching. The task sequences consist of compositional sequences of predefined behavior primitives shown in Figure 9 (a).

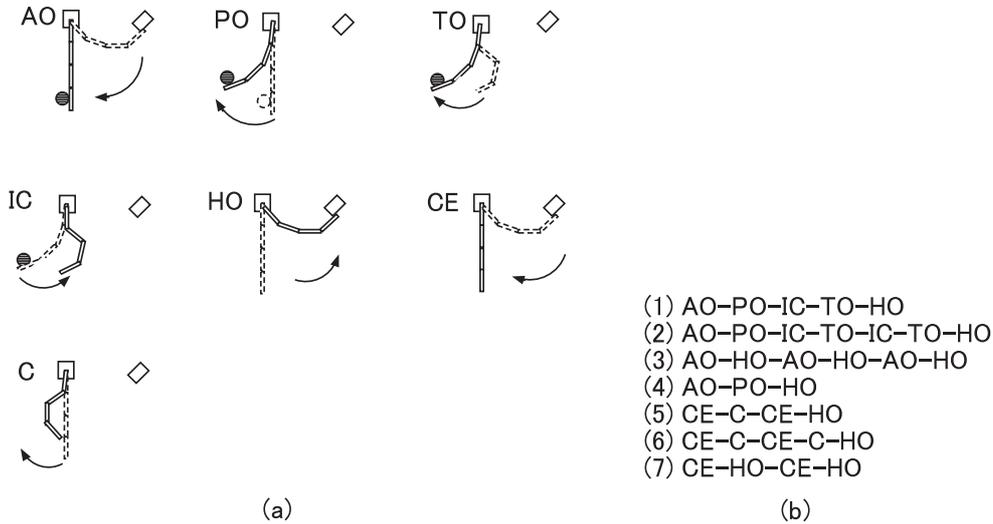


Figure 9: (a) Seven behavior primitives to be acquired. (b) The seven target task sequences which are composed of the seven primitives.

The seven primitives are AO: approach object in the center, PO: push object from the center to the left-hand side, TO: touch object, IC: perform inverse C shape, HO: return to home position, CE: move to the center, and C: perform C shape. Seven task sequences as shown in Figure 9 (b) are used for training of the network. These sequences contain branching structures such as – AO: approach object can be followed by either PO: push object or HO: return to home. Although there is no prior knowledge of the primitives in the network, the experimenter demonstrates them to the robot in terms of their combinational sequences. As a point of note, there are no explicit cues

for segmentation of the primitives in the continuous sensorimotor flow which the robot learns. The network must discover how to segment the flow of the ongoing task by attempting to decompose the sensorimotor flow into a sequence of segments (behavior primitives) which are reusable in other sequences to be learned.

6.2 Experiment

We begin by describing how the sensorimotor flow is segmented by self-organizing behavior primitives in the lower level by examining the PB activation sequences after the learning. We then present the results of additional experiments on real-time plan modulation utilizing the trained network for the purpose of examining dynamic characteristics of interplay between regression and prediction during actual behavior generation.

(A) Segmentation after learning

Figure 10 (a1-3) shows how the PB is activated during learning for three representative training sequences. The plots in the top row of this figure show the activation of four PB units as a function of the time step; the activation values from 0.0 to 1.0 are represented using a gray scale from white to black, respectively. The plots in the second and third rows represent the temporal profile of motor and sensor values for each training sequence. The vertical dotted lines indicate the occurrence of segmentation when the behavior sequence switches from one primitive to another in generating the training sequence. The capital letters associated with each segment denote the abbreviation of the corresponding primitive behavior. In this figure, observe that the switching of bit patterns in the PB takes place mostly in synchronization with the segmentation points known from the training sequences, although some segments are fragmented. Observe also that the bit patterns in the PB correspond uniquely to primitive behaviors in a one-to-one relationship in most cases.

(B) Online plan modulation

We examined how action plans can be dynamically adapted due to changes in the environment by going through bottom-up and top-down interactions. In this experiment, the robot learns to generate two different behavior sequences depending on the position of the object which is perceived as one of the sensory inputs. When the object is perceived in the center of the task space, the robot must perform task-A in which the arm repeatedly approaches the object and returns to the home position. This task can be represented as a primitive sequence of AO (approach object) followed by HO (return home). When the object is perceived on the left-hand side of the task space,

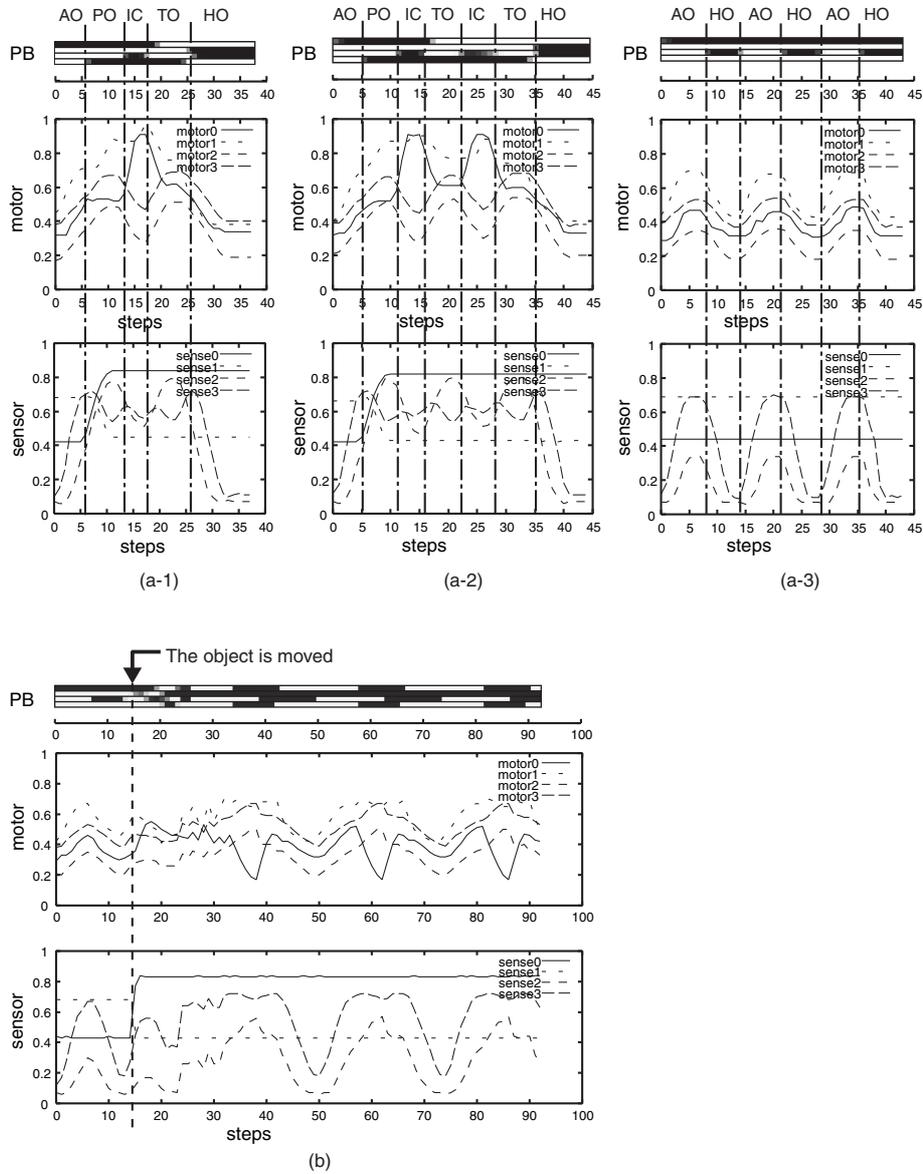


Figure 10: (a1-3) are three representative sequences after their training where the temporal profiles of the PB are plotted in the top row, the motor outputs in the second row, and the sensor inputs in the third row. The capital letters associated with each segment denote the abbreviation of the corresponding primitive behavior. (b) The profile of dynamic switching from task-A (before 13th step) to task-B (after 33rd step), where the dotted line indicates the moment when the object is moved from the center to the left-hand side.

the robot must perform task-B in which the arm repeats a primitive sequence of CE (move to center), C (make a C-shape), and then HO (return home).

The learning of these two tasks was attempted by training only the higher level while utilizing the behavior primitives learned in the lower level in the previous experiment. The lower level was not re-trained in this new experiment. After the learning was converged, it was shown that the robot could generate either of the behavior sequences correctly depending on the position of the object, that is, in the center or to the left-hand side. When the robot starts to move its arm, the PB vector in the initial event step is inversely computed to account for the sensory inputs representing the current position of the object as well as the current position of the hand tip. The two different situations concerning the object position generate two different PB vectors in the initial event step, which are followed by the corresponding PB sequences learned in the higher level. However, our current interest is to examine what would happen if the position of the object were to be switched from the center to the left-hand side in the middle of task execution. The question is that whether the behavior plans as well their executions can be dynamically adapted to the sudden situational changes through the process of the bottom-up and the top-down interactions.

The experiment showed that the task-A behavior pattern in the initial period is generated stably by following the mental plan. When the object is moved from the center to the left-hand side, there is an error in the anticipation of the sensory inputs, making the behavior pattern unstable and deviate from the task-A. However, this period of instability is resolved and the task-B behavior pattern initiates and it continues stably. Figure 10 (b) shows the temporal profile of the behavior generated. The vertical dotted line denotes the moment when the object is moved from the center to the left-hand side of the task space. Observe that it takes 20 steps until the task-B behavior pattern is initiated after the object is moved to the left-hand side. Observe also that the PB, the motor outputs, and the sensory inputs are distorted with fragmentation of the chunks during this transition period. Of particular note is the fact that the PB vector bit patterns are severely distorted.

A visualization of the fluctuation in video format is available as supplementary information at “<http://www.bdc.brain.riken.go.jp/tani/realtime>”. In this video, the fluctuation starts with the stepwise change in the sensory inputs which corresponds to the position change of the object introduced by the experimenter. It can be seen that the future plan image alternates every second between two possible temporal patterns of task-A and task-B. The PB vector also fluctuates not only for the current and future steps, but also for the immediate past steps due to the regression. The fluctuation is

initiated due to the gap generated between the top-down prediction of the PB values and their bottom-up estimation through the regression based on the sensory inputs. Repetition of this experiment 5 times with the same arbitration parameter k revealed the profiles of the transient motor patterns to be diverse, with the transient period ranging from 10 steps to more than 30 steps. Because of the unpredictable diversity in motor behaviors that occurred during this period, on one occasion the robot mistakenly knocked the object, making it fall from the table, at which time the experiment was terminated. It was, however, noted that such critical state while maximizing behavioral diversity appears only in a limited range of parameter k (see details in (Tani, 2003)).

6.3 An account of 'self-referential selves'

Let us now examine possible correspondences of the experimental results to the phenomenology of self-referential selves by firstly revisiting Husserl's thoughts on time perception. Husserl (Husserl, 1964) introduced the concepts of "retention" and "protension" to account for the subjective experience of "nowness". He explained the idea using an example of hearing the sound phrase "Do Mi So". When the note "Mi" is heard, we would still perceive a lingering impression of "Do", while at the same time anticipating hearing the next note of "So". The former is retention and the latter protension. These phenomena may also appear in physical behavior generation. For example, when our hands approach target objects from a resting position, we may have an impression of from where our hands have started just immediately before as well as a subtle expectation of touching the target object soon after. The terms retention and protension are used to designate the experienced sense of the immediate past and the immediate future, respectively. They should be considered as a part of automatic processing and thus unable to be controlled consciously. Husserl considered that the subjective experience of "nowness" is extended to include fringes in the experienced sense of both the past and the future in terms of retention and protension. Therefore, this "nowness" is not a point in physical linear time but has duration. The operation of the RNN shown in the previous experiments can be related to these phenomena of retention and protension. The reader is reminded that the prediction of the RNN is not made by an explicit logical computation, but rather proceeds as an autonomous dynamics while retaining the past context internally. This context-dependent forward dynamics of the RNN could be one possible realization of the phenomena of retention and protension (Tani, 2004).

If we understand Husserl's notion of "nowness" in terms of retention and protension,

the following question arises. Where is “nowness” bounded? Husserl seems to consider that the immediate past does not belong to a *representational* conscious memory, but merely to an impression. Yet how could the immediate past, experienced just as an impression, slip into the distant past which can be retrieved through a conscious memory operation (Varela, 1999)? What kind of mechanism qualitatively changes an experience from just an impression to a consciously retrievable objective event? Furthermore, Husserl’s goal was to explain the emergence of the objective time level from the pre-empirical level of retention and protension (Husserl, 1964). Husserl seems to consider that the sense of objective time would emerge as a natural consequence of organizing each experience into one consistent linear sequence. But, what are the underline mechanisms for this?

The idea of a segmenting flow of experience into objectified reusable units or primitives could be the key to answering these questions. Our main idea is that “nowness” can be bounded where the smooth flow of experience is truncated by conflicts in terms of prediction error. The sequential notes of “Do Mi So” constitute a chunk within which a perfect coherence is organized in the coupling between the neural dynamics of anticipation and the sound sensation flow. The same can be assumed for behavior generation such that, inside of the behavior chunk of “approach an object for grasping”, usually matching between anticipated sensation and actual one runs smoothly and automatically without giving rise to conflicts or ‘consciousness’. However, when we hear the second phrase “Re Fa La” after “Do Mi So”, a temporal incoherence emerges in the transition between the two phrases since this second phrase is not necessarily predictable from the first one. Similarly, after grasping the object, the chunk is truncated because of a possible immediate multiple branching proceeding to, for example, lift it up, pull it over, or throw it.

Let us consider this in greater depth by examining what happens in the hierarchical RNNPB. In the learning process, the moment of unpredictability arises when the target training sensorimotor flow encounters a branching where the RNNPB in the lower level cannot tell which way the flow proceeds. Remember that the training sequences contain branchings since they are designed to preserve compositionality. With this unpredictability, the PB vector is shifted from one value to another by the error generated. At this very moment, it can be interpreted that the automatic continuous sensorimotor flow experienced is segmented by the error associated with ‘consciousness’. Moreover, it can be said that *subjective* experience of the sensorimotor flow is *objectified* by means of the PB vector which can be manipulatable in the higher level. In the course of the learning process, the higher level becomes able to describe its own

experiences in linear sequences of events in terms of the PB vector sequences. In regards to memory retrieval, however, the sensorimotor flow can be reconstructed only in an abstract way since the original flow is now represented by combining a set of behavior primitives. Although such ways of reconstruction could provide compositionality as well as generalization in representing the sensorimotor flow experienced, the uniqueness of each instance of original experience might be lost by abstraction. Consequently, it is presumed that the sense of objective time may appear when the experience of the sensorimotor flow is reconstructed in a compositional form while losing its exactness.

In the previous section, it was conjectured that possible constructs for self-referential selves might be constituted via objectification of subjective experiences. Indeed, it was just shown that such objectification of sensorimotor flow experiences can be achieved through self-organization of the neuronal dynamics in the learning process. However, such constructs obtained may not yet account for genuine structures of self-referential selves since they were just constituted in a static way, along a one-directional bottom-up pathway. Incidentally, our previous experiment regarding the online plan modulation suggested that the sequencing of primitives in the higher level can become susceptible to unexpected perturbations, such as when an object is suddenly moved. Such perturbation could initiate complex situations. One aspect of complexity might arise from the online nature of behavior generation. As has been observed, if the top-down expectation of the PB values conflict with the those from the bottom-up regression based on the current experience, the PB vector can be fragmented. Even during this fragmentation the robot continues to generate behaviors, but in an abnormal manner due to distortion of the PB vector. The regression of this sort of abnormal experience causes further modulation of the current PB vector in a recursive way. During this iteration within the causal loop, the entire system may face intrinsic criticality, from which the observed diversity of behaviors might originate. We suggest that genuine constructs of self-referential selves could finally appear at such criticality accompanied by conflictive interactions with circular causality between the top-down subjective mind and the bottom-up reality.

7 Discussion

7.1 Minimal self, social self, and self-referential self

The current paper attempts to elucidate the dynamic characteristics of the autonomy of selves by reviewing a series of synthetic robotics experiments conducted by the author's

group.

In the neuro-robotics modeling of Experiment-1, mutual interaction between the bottom-up pathway of landmark perception and the top-down pathway of its prediction is arbitrated by internal parameters which are adapted by utilizing the prediction error. The experimental results revealed that the entire system dynamics proceeds with intermittent shifting between the coherent phase with good predictability and the incoherent phase with poor predictability through the incremental learning process. By referring to Heidegger's hammer example, it was postulated that the 'minimal self' becomes consciously aware by the gap generated between top-down anticipation and bottom-up reality in the incoherent period. It was also suggested that the entire system dynamics tends to proceed toward a certain critical state in which a large range of fluctuations could take place; a mechanism which might be analogous to SOC (Bak et al., 1987).

Experiment-2 explored characteristics of 'selves' in the social context by conducting an imitation game between a humanoid robot controlled by the RNNPB and human subjects. The RNNPB is characterized by its simultaneous processes of prediction of future and regression of past. In the middle of the mutual imitation game, spontaneous shifts were frequently observed between the states of coherence and incoherence accompanying turn taking phenomena in the interactions between the two sides. It was suggested that such complexity may appear at a certain critical period in the course of developmental learning processes by human subjects when an adequate balance between predictability and unpredictability is achieved. It was speculated that human subjects may perceive the autonomy of 'selves' for robots when they participate in interactive dynamics with criticality.

Experiment-3 addressed the problem of self-referential selves. Here, the RNNPB model was further extended with hierarchy. In the learning experiment with an arm robot manipulating an object, the continuous sensorimotor flow experienced was segmented into a sequence of reusable behavior primitives by accompanying stepwise shifts in the PB vector caused by the prediction error. Then, the higher level RNN learned to predict the sequences of behavior primitives in terms of shifts in the PB vector. The observed phenomena in the experiments could be interpreted as the process of achieving "self-referential selves". This is because the subjective experience of sensorimotor flow is objectified into reusable units which are manipulable in the higher level. Consequently, when the original experiences of sensorimotor flows are reconstructed with compositional structures, they become consciously describable objects rather than merely impressions of the original experiences. This might account for

how self-referential selves can be constituted.

It was, however, argued further that the genuine forms of self-referential selves should be constituted in causal circular interactions between the bottom-up pathway and the top-down pathway. In order to preserve rich interactions between the two, the entire system dynamics should be sustained at certain critical regions, which were depicted in our experiments of online plan modulations. It turns out that self-referential selves can evolve in diverse trajectories with their intrinsic autonomy characterized by such critical dynamics.

These experimental results suggest that although all three types of selves differ from each other, they also share a similar condition of criticality that emerges in dynamic interactions between the bottom-up and top-down processes.

7.2 Autonomy of selves

In this subsection, I would like to examine further how the autonomy of selves is constituted by exploring possible correspondences between phenomenology and the dynamical systems approach used in some synthetic modeling studies. For this purpose we firstly revisit Husserl's thoughts on temporality and Varela's interpretations of its dynamical systems (Varela, 1999).

One of the most essential but difficult considerations made by Husserl on temporality is so-called 'double intentionality' (Husserl, 1964). He attempted to provide accounts for the paradoxical characteristics of temporality of being both static (stable) and flowing (dynamic) by using the respective notions of transversal and longitudinal intentionality. Varela (Varela, 1999) interprets transversal intentionality as the static constitution of 'trajectories' appearing as event sequences relating the past, present and future in the object time level. He regards longitudinal intentionality as the dynamic property of self-motion, immanence or, said more simply, the continuously changing dynamic structure of constituting temporality. He argued that these two intentionalities are interdependent in the sense that dynamical structures generate trajectories, which in turn modulates the dynamical structures, and vice versa (Varela, 1999). He speculated that this sort of mutual bootstrap between trajectories and dynamic structures could lead to complex dynamics characterized as the edge of chaos (Crutchfield, 1989) or self-organized criticality (Bak et al., 1987).

Iizuka and Ikegami (Iizuka & Ikegami, 2004) reported an interesting simulation study which is relevant to our work. In their study, two mobile agents with RNNs were evolved with a genetic algorithm (Holland & Reitman, 1978; Nolfi & Floreano,

2000) to develop their interactive behaviors. They set an intriguing fitness function that maximizes the product of periods of one agent following and being followed by the other agent. Interestingly, this scheme resulted in the emergence of turn-taking behaviors in which the situations of following and being followed are spontaneously switched. It should be emphasized that the autonomy observed in the behaviors of these coupled agent are immanent because it self-manifests due to the evolutionary pressure that poses conflictive requirements for the agents to establish complementary actions. The micro-slip phenomenon modelled by Ogai and Ikegami (Ogai & Ikegami, 2008), again using an evolutionary scheme, might be explainable in the similar way. When agents are required to conduct two alternative actions with equal probabilities, the evolved behaviors of the agents show wondering trajectories between the two possible actions. Such conflictive situations imposed on the agents' actions lead to self-organization of chaos in the sensorimotor coupled neurodynamics.

Although the abovementioned studies nicely show an immanent property in the autonomy of agents in synthetic ways, this may not yet be for autonomy of selves and consciousness since it does not account for the issue of *temporalization*. As postulated by Husserl, selves and consciousness should be considered as inseparable from the perception of temporality. On this topics, Varela (Varela, 1999) wrote that:

Consciousness does not contain time as a constituted psychological category. Instead, temporal consciousness *itself* constitutes an ultimate substrate of consciousness where no further reduction can be accomplished.

It can be said that conscious selves should be immanent properties in *temporalization*. The flow of experience from now to the past is objectified in the bottom-up regression, while the thoughts objectified are reflected to ongoing experiences in the top-down anticipation from now to the future by accompanying enaction (see Figure 11.) At the very moment when these two streams mingle together through dense interactions, no separation can be made between the subjective mind and the external reality where unity of the temporal consciousness of selves could appear diversely. Our robotics experiments introduced above could provide a synthetic analysis of the parts of such dynamics. The dynamics tends to become critical under certain conditions due to the generation of conflicts between the regression to “represent” own experiences in the past objectively and future anticipation of them. It should be understood that “represent” signifies a fundamental intentionality in phenomenology which directs toward the objective time level starting from the absolute flow level through pre-empirical level, as postulated by Husserl. Further, “represent” is thought to be betrayed by new ex-

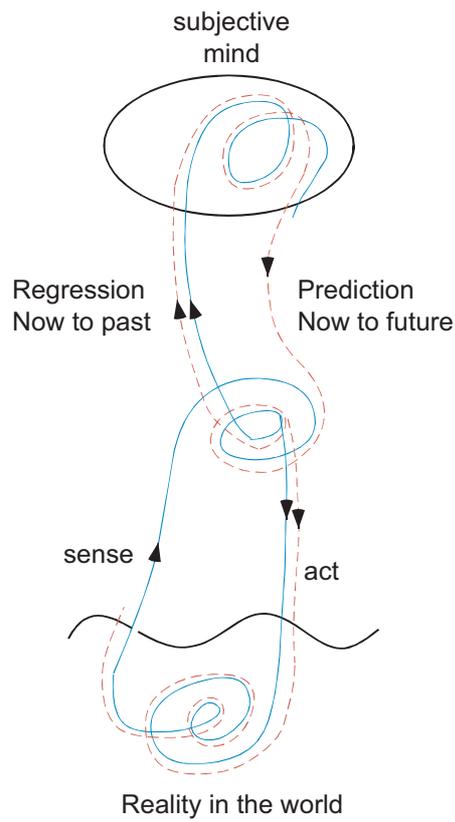


Figure 11: The two pathways of the top-down subjective mind of anticipating future and the bottom-up reality by regressing the past experience mingle together through their dense interactions.

periences and thus to be modified consistently. At the maximum of such fluctuations, conscious selves could emerge. I would suggest that everyday phenomena we encounter can be characterized as a never-ending race between experience and their objectification. Now, it is clear why I cannot find autonomy of 'selves' in behavior-based robots – they have no “representation” (Brooks, 1991) that make themselves struggle.

7.3 Criticality and authentic being

Before closing the current discussion, I would like to address one more topic which could help our understanding of the criticality discussed in the current paper – by viewing it from Heidegger's thoughts on existence (Heidegger, 1962), although it is purely speculative. Heidegger seems to consider that existentiality referred to as “being” or *dasein* is the final irreducible element in phenomenology. Therefore, it was suggested that thoughts on consciousness or selves should start from understanding “being” as a primordial phenomenon. Firstly, Heidegger considered that all things can “exist” in relational structures among others, like a hammer exists as a tool for a carpenter. Similarly, man can “exist” in a coherent relationship with his neighbors, for example, by engaging daily in “idle talk”. Heidegger considered that the *dasein* in this mode is simply an inauthentic being at the most, for man lives his daily life only in the immediate present, vaguely anticipating the future and mostly forgetting the past.

However, Heidegger also thought that this mode of *dasein* can be altered to that of an authentic one when man thinks positively about the possibility of his death, which can occur at any moment and not necessarily so very long in the future. Death is an absolutely special event because it is an ultimately individual event which cannot be shared with others. Although death can be regarded as an absolute impossibility of *dasein* which cannot be related to any other kinds of *dasein*, Heidegger considered that it can render the possibility of authentic *dasein* heading toward its absolute impossibility. When man looks ahead toward death as his ownmost possibility (the possibility of the impossibility of any existence at all) with anticipatory resoluteness (*vorlaufende Entschlossenheit*), he also needs to look back at his past with regression (*wiederholung*) in order to identify himself. Here, Heidegger's brilliant notion is that the present is “born” via dynamic interplay between *Zukunft* – looking ahead future for possibility and *Gewesenheit* – regressing past for reflection where authentic being is finally rendered. Consequently, Heidegger regards temporality as the ground for authentic being.

Although it might be regarded as absurd by the reader to relate our robots which

are technically regardless of their deaths to Heidegger’s thoughts, still we might find some correspondences by using some metaphor. Heidegger’s notion of authentic and inauthentic beings could account for the qualitative differences observed between our robots as reviewed in the current paper and conventional behavior-based robots. As described previously, the emergence of critical dynamics observed in our experiments can be accounted for by interactions between lookahead prediction of future and regression of past, which are parallel to Heidegger’s notion of the dynamic interplay between *Zukunft* and *Gewesenheit*. As the resultant behavioral trajectories during each task for the robot become diverse while exhibiting near power-law profiles which cannot be characterized by taking their average, their “being” might be said to be authentic metaphorically. On the other hand, conventional behavior-based robots simply trying to maintain a “nice” coherence with their environment cannot show any nontrivial behaviors because there is no intrinsic mechanism to drive them to criticality. Such robots might be regarded as “inauthentic beings”. An “authentic being” should be a primordial phenomenon that can render the genuine autonomy of selves and this should constitute the ultimate identity required for agency (Barandiaran et al., 2009) which I addressed in the introduction of the current paper.

References

- Arbib, M. (1981). Perceptual structures and distributed motor control. In *Handbook of Physiology: The Nervous System, II. Motor Control* (pp. 1448–1480). Cambridge, MA: MIT Press.
- Auvray, M., Lenaya, C., & Stewart, J. (2009). Perceptual interactions in a minimalist virtual environment. *New Ideas in Psychology, 27*, 32–47.
- Bak, P. (1996). *How nature works: The science of self-organized criticality*. Copernicus Books.
- Bak, P., Tang, C., & Wiesenfeld, K. (1987). Self-organized criticality: an explanation of the $1/f$ noise. *Phys Rev Lett, 59*, 381–384.
- Barandiaran, X., Paolo, E., & Rohde, M. (2009). Defining agency: individuality, normativity, asymmetry and spatio-temporality in action. *Adaptive Behavior, ???, ??-??* (current issue)
- Beebe, B., & Lachmann, F. (1988). Infant observation: The contribution of mother-infant mutual influence to the origins of self- and object-representations. *Psychoanalytic Psychology, 5*, 305–337.
- Beer, R. (1995). A dynamical systems perspective on agent-environment interaction. *Artificial Intelligence, 72*(1), 173–215.
- Blanke, O., & Metzinger, T. (2009). Full-body illusions and minimal phenomenal selfhood. *Trends in Cognitive Sciences, 13*, 7–13.
- Brooks, R. (1991). Intelligence without representation. *Artif. Intell., 47*, 139–159.
- Butz, M. (2008). How and why the brain lays the foundations for a conscious self. *Constructivist Foundations, 4*(1), 1–14.
- Crutchfield, J. (1989). Inferring statistical complexity. *Phys Rev Lett, 63*, 105–108.
- Gallagher, S. (2000). Philosophical conceptions of the self: implications for cognitive science. *Trends. CogSci., 4*(1), 14–21.
- Haruno, M., Wolpert, D., & Kawato, M. (2003). Hierarchical mosaic for movement generation. *International Congress Series, 1250*, 575–590.

- Heidegger, M. (1962). *Being and Time*. New York: Harper and Row.
- Holland, J., & Reitman, J. (1978). Cognitive systems based on adaptive algorithms. In D. Watermann & F. Hayes-Roth (Eds.), *Pattern directed inference systems*. New York: Academic Press.
- Holland, O., & Goodman, R. (2003). Robots with internal models a route to machine consciousness? *Journal of Consciousness Studies*, 10, 77–109.
- Hopfield, J., & Tank, D. (1985). Neural computation of decision in optimization problems. *Biological Cybernetics*, 52, 141–152.
- Horswill, I. (1993). Polly: A vision-based artificial agent. In *Proc. of AAAI-93* (pp. 824–829). MIT Press.
- Hume, D. (1975). *A treatise of human nature*. Clarendon Press.
- Husserl, E. (1964). *The phenomenology of internal time consciousness, trans. J.S. Churchill*. Bloomington, IN: Indiana University Press.
- Iizuka, H., & Ikegami, T. (2004). Adaptability and diversity in simulated turn-taking behavior. *Artificial Life*, 10(4), 361–378.
- Iizuka, H., & Paolo, E. (2007). Minimal agency detection of embodied agents. In *LNCS, Proc. ECAL2007* (Vol. 4648, pp. 485–494). Heidelberg: Springer.
- Ito, M., & Tani, J. (2004). On-line imitative interaction with a humanoid robot using a dynamic neural network model of a mirror system. *Adaptive Behavior*, 12(2), 93–114.
- James, W. (1890). *The principles of psychology*. Dover Publ. (reprinted 1950).
- Jordan, M. (1986). Attractor dynamics and parallelism in a connectionist sequential machine. In *Proc. of 8th annual conference of cognitive science society* (pp. 531–546). Hillsdale, NJ: Erlbaum.
- Kohonen, T. (2001). *Self-organizing maps*. Springer.
- Maes, P. (1991). *Designing Autonomous Agents: Theory and Practice from Biology to Engineering and Back*. Cambridge, MA: MIT Press.
- Martius, G., Nolfi, S., & Herrmann, J. (2008). Emergence of interaction among adaptive agents. In *LNCS Proc. of SAB2008* (Vol. 5040, pp. 457–466). Springer.

- Maturana, H., & Varela, F. (1980). *Autopoiesis and cognition: the realization of the living*. Boston: D. Riedel Publishing.
- Nolfi, S., & Floreano, D. (2000). *Evolutionary robotics: The biology, intelligence, and technology of self-organizing machines*. MA: MIT Press/Bradford Books.
- Ogai, Y., & Ikegami, T. (2008). Microslip as a simulated artificial mind. *Adaptive Behavior*, 16(2-3), 129–147.
- Paine, R., & Tani, J. (2005). How hierarchical control self-organizes in artificial adaptive systems. *Adaptive Behavior*, 13(3), 211–225.
- Reed, E., & Schoenherr, D. (1992). *The neuropathology of everyday life: On the nature and significance of microslips in everyday activities*. (unpublished manuscript)
- Rizzolatti, G., Fadiga, L., Galless, V., & Fogassi, L. (1996). Premotor cortex and the recognition of motor actions. *Cognitive Brain Research*, 3, 131–141.
- Rumelhart, D., Hinton, G., & Williams, R. (1986). Learning internal representations by error propagation. In D. Rumelhart & J. McClelland (Eds.), *Parallel distributed processing* (pp. 318–362). Cambridge, MA: MIT Press.
- Schoner, S., & Kelso, S. (1988). Dynamic Pattern Generation in Behavioral and Neural Systems. *Science*, 239, 1513–1519.
- Strawson, G. (1997). The self. *Journal of Consciousness Studies*, 4-5/6, 405–428.
- Tani, J. (1996). Model-Based Learning for Mobile Robot Navigation from the Dynamical Systems Perspective. *IEEE Trans. on SMC (B)*, 26(3), 421–436.
- Tani, J. (1998). An interpretation of the "self" from the dynamical systems perspective: a constructivist approach. *Journal of Consciousness Studies*, 5(5-6), 516–42.
- Tani, J. (2003). Learning to generate articulated behavior through the bottom-up and the top-down interaction process. *Neural Networks*, 16, 11–23.
- Tani, J. (2004). The Dynamical Systems Accounts for Phenomenology of Immanent Time: An Interpretation by Revisiting a Robotics Synthetic Study. *Journal of Consciousness Studies*, 11(9), 5–24.

- Tani, J., Ito, M., & Sugita, Y. (2004). Self-organization of distributedly represented multiple behavior schemata in a mirror system: reviews of robot experiments using RNNPB. *Neural Networks*, *17*, 1273–1289.
- Tani, J., & Nolfi, S. (1999). Learning to perceive the world as articulated: an approach for hierarchical learning in sensory-motor systems. *Neural Networks*, *12*, 1131–1141.
- Trevarthen, C. (1993). *The self born in intersubjectivity: The psychology of an infant communicating*. Cambridge: University Press.
- Turing, A. (1950). Computing machinery and intelligence. *Mind* *LIX*, *236*, 433–460.
- Varela, F. (1996). Neurophenomenology: a methodological remedy for the hard problem. *J. of Conscious Studies*, *3*(4), 330–350.
- Varela, F. (1999). Present-Time Consciousness. *Journal of Consciousness Studies*, *6*(2-3), 111–140.
- Wolpert, D., Doya, K., & Kawato, M. (2003). A unifying computational framework for motor control and social interaction. *Philosophical Transactions of the Royal Society B*, *358*(1431), 593–602.
- Yamashita, Y., & Tani, J. (2008). Emergence of functional hierarchy in a multiple timescale neural network model: a humanoid robot experiment. *PLoS Computational Biology*, *4*, e1000220.
- Ziemke, T., Jirnhed, D., & Hesslow, G. (2005). Internal Simulation of Perception: Minimal Neurorobotic Model. *Neurocomputing*, *68*, 85–104.