

From Biological to Synthetic Neurorobotics Approaches to Understanding the Structure Essential to Consciousness (Part 2)

Jun Tani¹

DEPARTMENT OF ELECTRICAL ENGINEERING, KOREAN ADVANCED
INSTITUTE OF SCIENCE AND TECHNOLOGY (KAIST), OKINAWA
INSTITUTE OF SCIENCE AND TECHNOLOGY (OIST),
TANI1216JP@GMAIL.COM

Jeff White

COMPUTATIONAL NEUROSYSTEM LABORATORY, KAIST

ABSTRACT

This paper reviews research in “predictive coding” that ultimately provides a platform for testing competing theses about specific dynamics inherent in consciousness embodied in both biological and artificial systems.

1 INTRODUCTION

We have been left with a big challenge, to articulate consciousness and also to prove it in an artificial agent against a biological standard. After introducing Boltuc’s h-consciousness in the last paper, we briefly reviewed some salient neurology in order to sketch less of a standard than a series of targets for artificial consciousness, “most-consciousness” and “myth-consciousness.” With these targets on the horizon, we began reviewing the research program pursued by Jun Tani and colleagues in the isolation of the formal dynamics essential to either. In this paper, we describe in detail Tani’s research program, in order to make the clearest case for artificial consciousness in these systems. In the next paper, the third in the series, we will return to Boltuc’s naturalistic non-reductionism in light of the neurorobotics models introduced (alongside some others), and evaluate them more completely.

1.1 PREDICTIVE CODING

In this section, we will review a research program into artificial consciousness that demonstrates the potential for computational experiments to isolate the formal dynamics of consciousness including the sense of time. Our focus is on the capacity for agents like human beings to project and to act towards possible futures by reflecting on the past. Studies in biological cognition have set out this capacity in terms of “predictive coding.”² With predictive coding, the results of actions—common “experience”—are integrated into an agent in terms of “prediction error.”

Prediction error informs the agent about how far from an intended target a prior action has led it, with the agent’s implicit aim being the minimization of this error signal. That said, minimization of error is not absolute. Optimizing for long-term ends may result in a relative detachment from the immediate perceptual reality, and conversely overt attention on immediate rewards may result in mounting error over the long run. Because predictive coding makes this form of future-oriented proactive agency based on effortful past regression possible within a mathematically embodied

agent, it offers a promising formal framework within which the relationship between the subjective mind and the objective world may be instantiated in an artificial agent.

Predictive coding is an important development in artificial consciousness research in two important ways. One, it provides a direct way to model subjective intention within the objective world. And two, it provides an equally direct way to project back the reality of the objective world as perceived by and as consequent on the actions of embodied and embedded cognitive agents.³ The result is a fully accessible dynamical mirror into the operations essential to consciousness in more complex systems, a promise that merely biological approaches to the study of consciousness cannot match. Tani was the first to successfully instantiate predictive coding in artificial agents, e.g., robots, in a deterministic domain, i.e., where intended outcomes are stable attractors.⁴ Alternatively, Friston explored Bayesian predictive coding in a probabilistic domain and generalized it under the name of the “free energy minimization principle” (FEMP).⁵

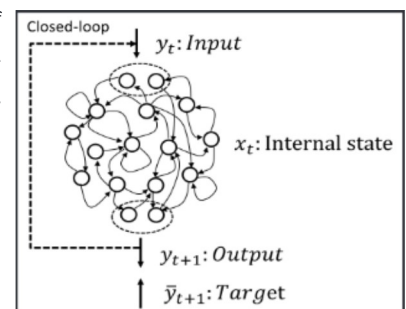
In the next section, we will briefly review a dynamic neural network model, the recurrent neural network (RNN),⁶ because it is a basic component of contemporary intelligent systems, and central to Tani’s deterministic dynamics which is the subject of the subsequent section. This review should serve as a primer on the dynamic system’s approach to embodied cognition. After reviewing Tani and colleagues’ formulation using RNN models, we will examine Bayesian predictive coding as formulated by Friston and colleagues.

2. PREDICTIVE CODING IN DYNAMIC NEURAL NETWORK MODELS

2.1 THE RNN MODEL

The essential characteristic of the RNN⁷ is that it can generate temporal sequence patterns as targets embedded in its internal dynamic structure. It “learns” to imitate exemplar sequence patterns, and when properly organized even to creatively compose its own⁸ by extracting underlying regularity. An example of an RNN is shown in Figure 1. This figure shows an RNN used in the predictive learning scheme to be described later (section 2.2).

Figure 1. RNN model of predictive learning with teaching target. The dotted line represents a closed-loop copy from the output to the input.



An RNN consists of a set of neural units including input units representing the input state, internal (context) units representing the internal state and output units representing the output state. These are variously interconnected by synaptic connectivity weights. These connections can be unidirectional, bidirectional, or recurrent. The time

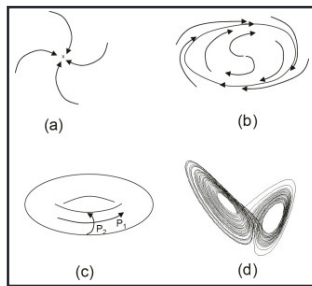
development of each neural unit output activation in discrete time can be written as:

$$u_i^{t+1} = \sum_j^N w_{ij} a_j^t + b^i \quad (1-a)$$

$$a_i^t = f(u_i^t) \quad (1-b)$$

Where u_i^t is the internal state of the i th neural unit at time step t , a_i^t is its output activation, w_{ij} is synaptic connection weight from the j th unit to the i th unit, b^i is the bias of the i th unit. $f()$ is a sigmoid function. Over time, the neural activation of the whole network can generate different types of dynamic attractor patterns depending on the synaptic weights adopted in the network. Figure 2 shows typical attractors including a fixed point attractor, a limit cycle attractor, a limit torus, and chaotic attractor.

Figure 2. Four different types of attractors. (a) fixed point attractor, (b) limit cycle attractor, (c) limit torus characterized by two periodicities P1 and P2 which form an irrational fraction, and (d) chaotic attractor.



The simplest attractor is a fixed point attractor in which all dynamic states converge to a point (Fig. 2 (a)). The second one is a limit cycle attractor (Fig. 2 (b)) in which the trajectory converges to a cyclic oscillation pattern with constant periodicity. The third one is a limit torus that appears when there is more than one frequency involved in the periodic trajectory of the system and two of these frequencies form an irrational fraction. In this case, the trajectory is no longer closed and it exhibits quasi-periodicity (Fig. 2(c)). The fourth one is a chaotic attractor in which the trajectory exhibits infinite periodicity and thereby forms fractal structures (Fig. 2 (d)).

These different types of attractor dynamics can account for the autonomous generation of different types of agent action patterns. For example, fixed point attractor dynamics account for a hand reaching movement, from an arbitrary hand posture to its end point, while limit cycle attractor dynamics account for a rhythmical hand waiving pattern with a certain periodicity, and chaotic attractor dynamics account for non-periodic, seemingly random movement.

RNNs can learn to generate such attractor dynamics through predictive learning. Each specific attractor pattern can be developed in an RNN by optimizing the synaptic weights and biases through a process of error minimization. In predictive learning, the network receives current time step perceptual input and outputs a prediction of the next time step (see Fig.1). Error is computed between the predicted output and the target (e.g., teaching exemplar), and synaptic weights and biases are updated in the direction of minimizing this error using error back-propagation through time (BPTT).⁹ After learning, the RNNs internal dynamics converge on a stable pattern, and the learned attractor can be generated from a given initial state through “closed-loop” (off-line) operation in which the predicted output of

the current time step is copied to the input of the next time step in a closed-loop (see the dotted line in Fig.1). This closed-loop operation corresponds to mental simulation, as will be described later sections.

An RNN can be regarded as a dynamical system with adaptive parameters including synaptic weights and biases which can be described in the following generalized form

$$x_{t+1} = F(x_t, w) \quad (2a)$$

$$y_{t+1} = G(x_{t+1}, W) \quad (2b)$$

In these expressions, x_t and y_t represent the current internal state and the output state, respectively, and w stands for the adaptive parameter. The internal state x_t is important, because it represents the current context or situation for the system as a whole and develops by means of an iterative learning process. The system can exhibit contextual information processing through which the output of the system reflects not merely the immediately perceived inputs but the context accumulated over past experiences of inputs. Formally speaking, this system embodies temporality, entrained according to patterns that extend beyond the immediate context and, as we shall see, reaches—even creatively, and inferentially—toward goal states.

The conventional RNN model can learn to generate only a single attractor pattern except special cases of developing multiple attractors. So, a natural question arises: How can the model be advanced such that it can learn to generate multiple attractor patterns, each specific to a different context? This question motivated an investigation into the possibility of applying the framework of predictive coding in the advancement of RNN models, as described next.

2.2 MIXRNNs AND RNNPB

Tani and colleagues investigated how a network model can retrieve and generate a particular sequential pattern from long-term memory of multiple patterns. Two versions of RNN models resulted, namely a mixture of RNN experts (MixRNNs)¹⁰ and a recurrent neural network with parametric bias (RNNPB).¹¹ MixRNNs use a local representation scheme, and RNNPBs use a distributed representation scheme, in order to learn to generate and to recognize sequences of primitive action patterns. Moreover, these movement patterns are temporal patterns requiring active self-entrainment through online information by another’s live and more-or-less similarly embodied example, recalling the mechanism of “mirror neurons.”¹² In this section, we look more closely at how the MixRNN and the RNNPB capture aspects of consciousness typically associated with more complex biological systems.

MixRNNs¹³ consist of sets of local RNNs internally associated through gates where the global output of the whole network model is computed as the weighted sum of the gate opening ratio for all local RNN outputs (see Figure 3).

During learning, local RNNs compete to predict the next perceptual state as the gate opens most for the local RNN with the least prediction error. Because the learning rate of

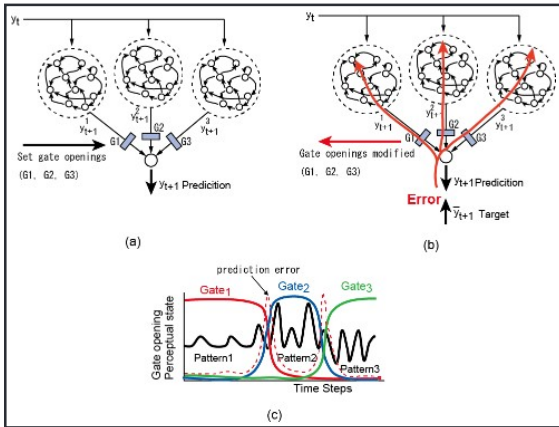


Figure 3. Description of MixRNNs. (a) Generation mode (b) recognition mode, and (c) segmentation of perceptual flow into a sequence of chunked sub-patterns by inferring gate openings.

each RNN is proportional to the gate opening ratio, the more that the gate of a particular RNN opens, the more this local RNN is able to learn the current perceptual sequence pattern. The goal of learning is to obtain optimal synaptic weights for all modular RNNs as well as optimal openings of the all gates at each time step, and by "optimal" we mean those which minimize the reconstruction error between the global output and the target output.¹⁴ Through a competitive learning process, i.e., error regression training with BPTT for the optimal gate opening sequence between RNNs, as well as for optimal synaptic weights in all local RNNs, each local RNN becomes an expert for a particular perceptual sequence pattern. Intuitively then, MixRNNs can learn a set of frequently apparent primitive patterns with each consolidated in a corresponding local RNN simply through the iterative and collective experience of those patterns.

After learning, a MixRNN model can generate a particular intended perceptual sequence by opening the gate of the corresponding RNN expert (Fig. 3(a)). In this way, current gate openings represent the current top-down intention designating the pattern to be generated. Additionally, MixRNNs can recognize a given perceptual sequence pattern through competition between local RNNs by reconstructing the target pattern with the least error by means of the error regression scheme optimizing the gate openings (Fig. 3 (b)), with synaptic weights fixed in this case). When error is minimal, a gate associated with a particular local RNN opens in a winner-take-all manner, and the target pattern is recognized as belonging to this expert RNN. In other words, the target pattern of the current perception can be recognized by means of reconstructing it in a particular local RNN with minimum error whereby the current gate opening states represent the inferred intention.

When the currently perceived sequential pattern changes, gate opening is shifted toward minimizing prediction error arising at this moment. An important point here is that the continuous perceptual flow is segmented into chunks by means of gate openings during these moments. Tani and Nolfi argue that this suddenly required effort for minimizing the error by inferring appropriate gate openings should accompany momentary consciousness.¹⁵ Next, we look at a further advance on RNNs in this direction, the recurrent neural network with parametric bias, or RNNPB.

The RNNPB¹⁶ is a single RNN model employing parametric bias (PB) units (see Figure 4).

PB represents the current intention as it projects a particular perceptual sequential pattern onto the external world, analogous to the gate dynamics in MixRNNs. PB does this job by playing the role of bifurcation parameter modulating the dynamical structure of the RNN.

In simple terms, an RNNPB learns to predict or generate a set of perceptual sequence patterns associated with corresponding PB vector values. During learning, the optimal synaptic weights for all different sequence patterns as well as the optimal PB vector value for each sequence pattern can be obtained. After learning, an RNNPB can generate a learned perceptual sequence pattern by adopting the PB with the corresponding vector value (Fig. 4(a)). It can also recognize a perceptual sequence pattern given as a target by inferring the optimal PB vector by way of which the

target sequence can be reconstructed and output with the minimum error (Fig. 4(b)). Fig. 4(c) shows how the continuous perceptual stream can be segmented into a sequence of prior-learned patterns in terms of attractor dynamics by tracking modulations in PB vector bifurcation parameters at each time step.

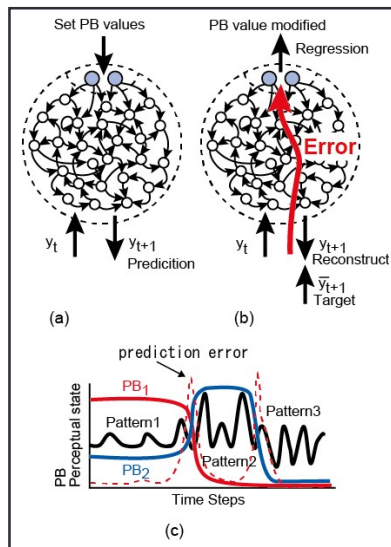


Figure 4. Description of RNNPB. (a) Generation mode, (b) recognition mode, and (c) segmentation of perceptual flow by PB vector into chunked patterns.

In the end, switching between chunks in the RNNPB is analogous to the segmentation mechanism employed in MixRNNs which use gates between local RNNs to recruit the appropriate expert or combination of experts given immediate perceptual reality. One limitation common to both is that each consists of a single level. But, when they are organized into a hierarchy, they can exhibit higher-order cognitive competencies such as creative compositionality.

Such an extension is the subject of the next section.

2.3 FUNCTIONAL HIERARCHY IN THE MTRNN

Both MixRNNs and the RNNPB have been developed into multiple-level architectures.¹⁷ The basic idea is that higher levels attempt to control lower levels by projecting control parameters (such as gate openings or PB vector modulations) onto lower levels based on current higher

order intention. And in turn, during normal operation the prediction error generated against the perceptual reality in the lower level is back-propagated to the higher level, where the control parameters for the lower level as well as the intention state in the higher level is adjusted in the direction of minimizing the error, i.e., by conforming to that state which would have resulted in least error.¹⁸

Tani and Nolfi demonstrate that hierarchically organized MixRNNs can learn to generate and recognize a set of sequential combinations of movement primitives in a simulated indoor robot navigation space.¹⁹ The analysis showed that a set of chunks related to movement primitives such as turning to the right/left at a corner, going straight along a corridor, and passing through a T-branch developed in local lower level RNNs, while different sequential combinations of these primitives developed in the higher-level RNNs, e.g., traveling through different rooms. When the simulated robot, for example, turns left at a corner from a straight corridor in a particular room, the continuous perceptual flow of its range sensor is segmented into the corresponding two movement primitives in the lower level. On the other hand, when it travels from a familiar room to another, segmentation related to the room transition can take a place in the higher level.

Tani achieved similar results in a real robotic arm with a similarly hierarchically organized RNNPB, which was able to deal with primitives and their sequential combinations during a simple object manipulation task. It is important to note that what begins as raw experience of the continuous perceptual flow becomes a manipulable object for the higher level after segmentation into chunks. Thus, the hierarchical structure adopted by Tani enables the objectification of perceptual experience, as will be described in greater detail later.²⁰

Building on this work in hierarchically organized RNNs, Yamashita and Tani²¹ demonstrated the learning of compositional action generation by a humanoid robot employing a novel multiple timescale RNN (MTRNN) (Figure 5). This MTRNN model uses multiple timescale constraints, with higher-level activity constrained by

slower timescale dynamics, and with lower level activity proceeding according to faster timescale dynamics. The basic idea is that higher-level information processing becomes more abstract as constrained by its slower dynamics, whereas lower level information processing is more sensitive to immediate details as constrained by faster dynamics.

The MTRNN shown in Fig. 5(a) consists of 3 subnetworks (slow, intermediate, fast dynamics networks; note that the numbers of levels can vary depending on application) each consisting of leaky integrator neural units that are assigned different time constants. The activation dynamics of a leaky integrator neuron can be described as a differential equation:

$$\tau \dot{x} = -x + \sum_i w_i x_i + I$$

$$\tau = 1 / (1 + \tau_0^2)$$

(3a)
(3b)

where τ represents the time constant of the neuron. When τ is set with a larger value, the time development of the neural activation becomes slower. With a smaller value, it becomes faster. Eq.3 is integrated over time using the difference method. The fast dynamics network in the lower level consists of two modular RNNs, one for predicting the proprioceptive state in the next step from current joint angle information, and the other for predicting low dimensional visual features in the next time step from current visual information.

During these humanoid robot learning experiments, the MTRNN was trained to generate a set of different visuo-proprioceptive trajectories corresponding to supervised targets by optimizing connectivity weights as well as the intention state corresponding to each trajectory. The intention state here is analogous to the PB value in the RNNPB, and corresponds with the initial states of neural units in the slow dynamics network of the MTRNN (see Fig. 5(a)). When learning begins, for each training sequence the initial state of the intention units is set to a small random value. The forward top-down dynamics initiated with this temporarily set initial state generates a predictive sequence. The error generated between the training sequence and the output sequence is back-propagated along the bottom-up path through the subnetworks with fast and moderate dynamics to the subnetwork with slow dynamics. This back-propagation is iterated backward through time steps via recurrent connections, whereby the connection weights within and between these subnetworks are modified in the direction of minimizing the error signal (at each time step). The error signal is also back-propagated through time steps to the initial state of the intention units, where these initial state values are modified.

Here, we see that learning proceeds through dense interactions between the top-down regeneration of training sequences and the bottom-up regression through these sequences by way of error signals, just as in the RNNPB. And as a result of this interaction, the robot learns a set of behavior primitive patterns such as reaching for an object, lifting the object up and down, or moving it left and right. These develop as distributed activation patterns in fast and intermediate dynamics networks while various control sequences for manipulating these primitive constructs

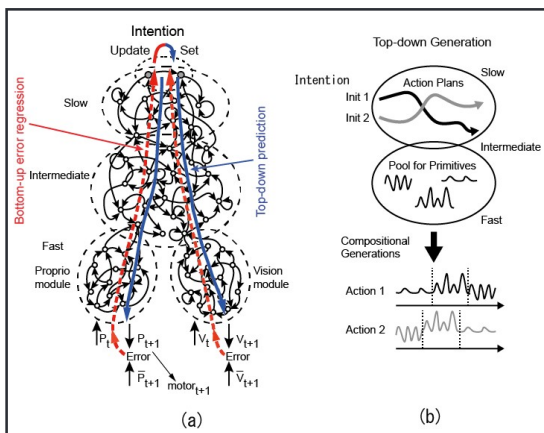


Figure 5. MTRNN model. (a) MTRNN architecture consisting of 3 levels, and (b) its top-down compositional generation of different intended actions.

develop in the slow dynamics network (according to its initial sensitivity characteristics, see Fig. 5 (b)).

What explains the success of these models in performing such complex cognitive tasks? In the MTRNN, neural activity output from the higher level plays the role of bifurcation parameter for the lower level, like the PB vector in the RNNPB. Building from this work, Yamashita and Tani concluded that the decomposition of complex visuo-proprioceptive sequences into sequences of reusable primitives is achieved within this functional hierarchy due to subnetwork timescale differences.²² Further experiments by Nishimoto and Tani and Arie and colleagues showed that MTRNNs can not only generate actual movements, but also diverse mental simulations of various intention states by performing closed-loop look ahead (so-called "off-line") prediction.²³ So, the question now becomes how to understand such functional hierarchies.

The development of functional hierarchies is captured in a well-known concept central to the study of complex adaptive systems, "downward causation," the causal relationship from global to local parts of a system.²⁴ A functional hierarchy develops by means of upward causation in terms of collective neural activity, both in forward activation dynamics and in error back-propagation. In the other direction, this development is subject to downward causation in terms of timescale difference, network topology, and environmental interaction. Note that these are strictly deterministic features of the system. Target conditions are determined. Current states are determined, and thereby optimal sequences of action are inferred. Next, we will look at an effort to articulate these temporal dynamics nondeterministically, in Friston's Bayesian predictive coding scheme formulated according to the free energy minimization principle.

3. THE FREE ENERGY MINIMIZATION PRINCIPLE

From the subjective perspective of an agent in the world, phenomena may be better described probabilistically than deterministically. Where upcoming anticipated optimal conditions are not pre-determined or perhaps even pre-determinable, the aforementioned models by Tani and colleagues can be extended into the probabilistic domain, as Friston has done.²⁵ Friston's main idea is to predict the next time step's perceptual states in terms both of their averages and their variances (or estimated accuracy). The average is a value arrived at according to prior instances, and actions undertaken on the basis of averages succeed best when deviations from the average are minimal. Variance, on the other hand, is a measure of the amount of difference between instances, and so can represent the accuracy of a prediction. Specifically, it can represent the estimated accuracy of a prediction, as a form of second-order prediction.²⁶

Now, the exact formula for representing this idea is derived from the principle of free energy minimization.²⁷ Negative free energy F can be computed by the addition of Gibbs energy G and Entropy E :

$$F = G + E \quad (4-a)$$

Then, F can be written in the following form.

$$F = -\sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) - \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t}))$$

Where $q(z)$ represents the prior distribution of the intention state, $P_{\theta}(x, z)$ represents joint probability distribution of observation x and the intention state Z parameterized by parameter θ . Then, negative free energy F can be transformed as:

$$\begin{aligned} F &= -\sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) - \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t})) \\ &= -\sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) - \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t})) \\ &= -\sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) - \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t})) + \sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) - \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t})) \quad (4-b) \end{aligned}$$

The last form obtained in (4-b) is equal to the lower bound, L which is well known in the machine learning field. The first term represents the likelihood of reconstructing X by parameter θ and the second term represents minus KL divergence between the prior probability distribution of the intention state and the posterior distribution after observation of X with parameter θ . It can be seen that maximizing the negative free energy is equal to maximization of the lower bound. This lower bound L can be rewritten as:

$$L = \sum_{s=1}^S \sum_{t=1}^T \ln(p(x_{s,t} | z_{s,t})) / 2 \ln(\frac{p(x_{s,t} | z_{s,t})}{q(z_{s,t})}) + \sum_{s=1}^S \sum_{t=1}^T \ln(q(z_{s,t})) / 2 \ln(\frac{p(x_{s,t} | z_{s,t})}{q(z_{s,t})}) \quad (4-c)$$

where $o_{s,t,i}$ is the i th dimension of the prediction output at time step t in the s th sequence, $\bar{o}_{s,t,i}$ is its teaching target, and $v_{s,t,i}$ is its estimated variance, $l_{s,t,i}$ is the i th dimension of the intention state for the s th sequence, and $\hat{\sigma}_s$ is its predefined deviation.

The generation, recognition, and learning of complex action sequences are possible through the maximization of negative free energy in the probabilistic domain just as the minimization of error performs similarly in the deterministic domain. According to the first term on the right-hand side of equation (4-c), the likelihood part can be maximized if variance is taken to be large even if the prediction square error is large. In this case, the agent has no reliable guide to anticipated future situations, so it simply relaxes any expectation of oncoming events. This would correspond with a reactive posture in a biological consciousness, for example. On the other hand, the likelihood might be small even though the prediction square error is small if the estimated variance is smaller than reality. In this case, an agent acts from intentions as if ends are predetermined, e.g., as if he has plotted all the necessary dimensions and their internal deviations so that action is facilitated and success presumed guaranteed. But, the agent ends up wrong about this, and suffers the correction. In human experience, having failed to adequately account for the world while having proceeded with laid plans in confidence is called "surprise." Similarly, according to Friston's free energy minimization principle (FEMP), the prediction square error divided by estimated variance represents the degree of surprise with interesting implications for inquiry into consciousness. For one thing, the measure of surprise may correlate with a measure of consciousness as the top-down accommodation of perceptual inputs at each time step.

According to the second term on the right-hand side of Eq. (4-c), the distance from the prior to the posterior can be minimized when the intention state of each sequence is distributed by following the Gaussian distribution with the predefined deviation δ_{is} . The recognition of the intention in FEMP is to infer the optimal probabilistic distribution of the intention state for a given target sequence, maximizing negative free energy rather than infer a single optimal value minimizing the prediction error as in the RNNPB. Instantiating such a process in a model dynamic system is subject of the next section.

3.1 THE STOCHASTIC MTRNN MODEL

Because the original FEMP by Friston²⁸ was not implemented in any trainable neural network models, it was not clear how maximizing negative free energy in Eq.(4-a) might lead to successful learning of internal predictive models extracted from perceptual sequence data experienced in reality. For this reason, Murata and colleagues proposed a novel dynamic neural network, referred to as the stochastic-MTRNN (S-MTRNN) model.²⁹ This model incorporates Friston's (2010) FEMP into the deterministic learning model described in the last section, the MTRNN. The S- extends the original as it learns to predict subsequent inputs taking into account not only their **means** but also their "variances," or range of anticipated values. This means that if some segments of input sequences are more variable than others, then the time-dependent variances over these periods become larger. On the other hand, if some parts are less fluctuated, their variances are smaller. In effect, then, the S-MTRNN predicts the predictability of its own prediction for each dimension of the input sequences in a time-dependent manner. When variances are estimated as zero, then the S-MTRNN becomes a deterministic dynamic system like the original MTRNN, i.e., it anticipates zero variance. Therefore, it can be said that—depending on context—S-MTRNNs can develop either deterministic or stochastic dynamics, at which point arises the notion of probability and so some valuation of possible future states accordingly.

The model operates by means of maximizing the negative free energy described in Eq. (4) in all phases of learning, recognizing, and generating perceptual sequence patterns. An important development in the current model is that $v_{s,i}$ as estimated variance in the likelihood part of Eq. (4-c) is changed to a time variable valuable $v_{s,i,t}$ because its estimate can change at each time step of a perceptual sequence. The likelihood part exists to minimize the square error divided by estimated variance at each step. This means that the prediction error at a particular time step is pressured to be minimized more strongly when its estimated variance is smaller. Otherwise, the pressure to minimize prediction error is less.

Another development is that the intention state $I_{S_{s,i}}$ in the part of KL divergence between the prior probability distribution of the intention state and the posterior distribution in Eq. (4-c) is now represented by the initial states of context units in all levels. The KL divergence part of Eq. (4-c) puts specific probabilistic distribution constraints on optimal initial states for all sequences with the parameter α_{is} . If α_{is} is set with a large value, the distribution of initial states becomes wide. Otherwise, it becomes tight. By maximizing the

negative free energy during the learning process, optimal connectivity weights for all teaching sequences, the probability distribution of the initial state for each teaching target sequence,³⁰ and the estimates of time-dependent variance for each sequence are obtained.

Figure 6 (a) shows the architecture of the S-MTRNN. The difference from the original MTRNN is that the S-MTRNN contains output units for predicting variances for all sensory dimensions at each time step.

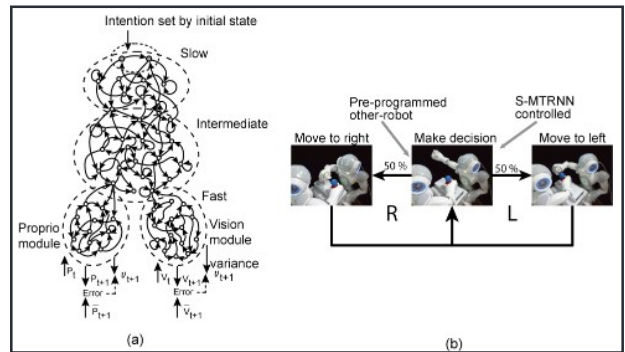


Figure 6. S-MTRNN model and the robotic experiment with the model. (a) The S-MTRNN contains additional output units for predicting variances for all sensory dimensions at each time step. (b) The "self-robot" learns to generate cooperative behaviors with the "other-robot."³¹

The next section reviews how the S-MTRNN performs in a particular robot task in the probabilistic domain.

3.2 LEARNING TO COOPERATE WITH OTHERS

A robotic experiment was conducted utilizing the S-MTRNN described in the preceding section (see Fig. 6).³² The objective of this robot experiment was to examine how one robot can generate "cooperative" behavior by adapting to another robot's behavior, even though its predictions occasionally fail. The experiment used two "NAO" humanoid robots. One NAO robot, the "self-robot," attempted to generate cooperative behaviors with the "other" NAO robot. The other-robot's behavior was pre-programmed. The self-robot was controlled by the S-MTRNN model.

During the experiment, the other-robot repeated movement patterns and the self-robot was tutored to generate corresponding "cooperative" behaviors as it perceived the other's object movements. In order to do this, the self-robot needed to proactively initiate its own arm movement before sensing the actual movement initiated by the other-robot. The self-robot acquired this cooperative behavior skill through direct tutoring from the experimenter.³³ The self-robot observed the other-robot perform sequences of five movements, moving a colored object either to the left or to the right in all possible combinations (2^5 sequences). Then, the self-robot was required to generate cooperative behaviors by simultaneously moving its arm in the same direction as the other-robot. As it generated movements and adjusted to the other-robot's movements, differences emerged in the dynamics involved in predicting as well as

generating behavior between the two conditions of the wide and narrow initial states. The following explains how we tested these results in greater detail.

During the test phase of the experiment, the S-MTRNN was trained with 2^5 visuo-proprioceptive (VP) sequences during the tutoring process. This training was repeated twice, once with a small value for σ and then again with a large value in order to generate a narrow initial state distribution (Narrow-IS) and a wide initial state distribution (Wide-IS), respectively. Other-robot object movement (either to the left or to the right) was randomly determined from amongst the same 2^5 sequences so that the self-robot (S-MTRNN) would be unable to predict next movements reliably.

After training, closed-loop generation of “mental” imagery was performed for both wide and narrow training cases (i.e., offline rehearsal). During closed-loop operation, Gaussian noise corresponding to the estimated variance at each step was applied to the feedback from the previous step prediction output, and was input to the current step prediction (see Figure 7 (a)).

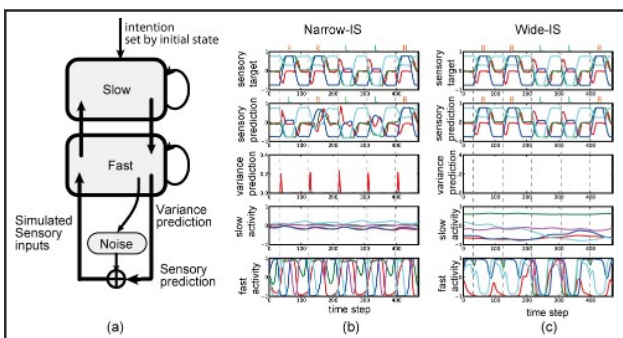


Figure 7. Generation of mental imagery via closed loop. (a) Closed-loop generation by S-MTRNN, generated sequences (b) in Narrow-IS case and (c) in Wide-IS case.³⁴

In this way, mental imagery increasingly fluctuates as uncertainty of a prediction, i.e., the estimated variance, increases. In the example pictured in Fig. 7, the initial states were set with the values obtained upon learning the “RRLLR” trial sequence as performed by the other-robot. Fig. 7 (b) and (c) illustrate mental imagery in terms of prediction of VP sequences associated with estimated variance and internal neural activities in the fast and the slow subnetworks as generated by the S-MTRNNs trained under both Narrow-IS or Wide-IS conditions. In the Narrow-IS case, diverse decision sequences were generated even though all trials began from the same initial state. As the figure shows, estimated variance sharply peaks at decision points, but remains almost zero at other time steps. This implies that during training the S-MTRNN develops action primitives for moving left or right as two distinct chunks, and employs a probabilistic switching mechanism at decision points.

On the other hand in the Wide-IS case, the same decision sequence was repeatedly generated for the same given initial state. Fig. 7 (c) shows that the VP sequence for “RRLLR” was generated which seemed to be mostly the

same as the target VP sequence. Here, it is important to note that the variance is estimated as almost zero for all steps including at decision points. This implies that mental imagery is generated as a deterministic predictive dynamics in the Wide-IS condition. Interestingly, for more than 20 branching instances before finally converging to cyclic branching, the robots’ “mental imagery” (predictive dynamics) of next-movements was generated pseudo-randomly by means of transient chaos that developed in the slow dynamics part in the model network. This result is analogous to that of Namikawa et al., where complete chaos (with a positive Lyapunov exponent) instead of transient chaos appeared in the neural dynamics of an MTRNN.³⁵

In the end, neural activity internal to the Narrow-IS and Wide-IS systems was quite different. In the Narrow-IS case, the neural activities in both the slow and the fast subnetworks showed the same values at all decision points. In the Wide-IS case, slow and fast neurons exhibited different activation patterns at each decision point through which the system was able to attempt to predict the subsequent move, left or right. There appears to be no such bias in activity at decision points in the Narrow-IS case, whereas there are top-down predictive biases imposed by specific top-level neural activation patterns at decision points in the Wide-IS case.

Let’s look more closely at how the self-robot interacted with the other robot using the network trained in these two conditions, Wide-IS and Narrow-IS. Starting with arbitrary initial states, S-MTRNN generated one-step predictions for subsequent VP states upon perceiving current visual states via open-loop generation, while the other robot randomly moved a colored object sequences of five (see Figure 8 (a)).

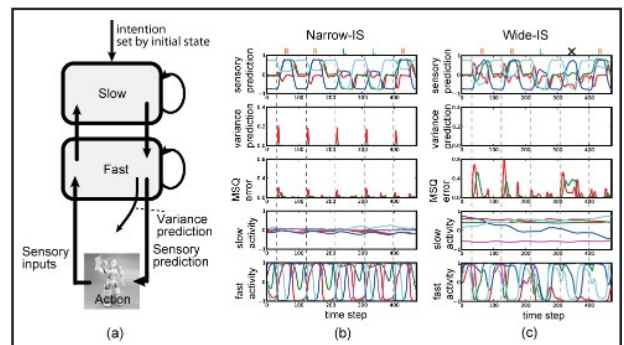


Figure 8. The results of the self-robot interacting with the other-robot by open-loop generation. (a) The scheme of the open-loop generation, (b) a sequence generated by the network trained with the Narrow-IS condition and (c) with the Wide-IS condition.³⁶

Fig. 8 (b) and (c) show the results of open loop processing with the self-robot reacting to the other-robot as it generated the “RRLR” sequence for the Narrow-IS and the Wide-IS cases, respectively. Here, we observe that one-step prediction of VP states in the Narrow-IS case is quite successful, generating only a small error at each decision point. In contrast, one-step prediction in the Wide-IS case is much worse. In fact, the prediction error is significantly large at many decision points. Interestingly, at this juncture of the trials, the movement of the self-robot became erratic.

For example, in the fourth decision as illustrated in Fig. 8 (c), the self-robot moved its arm in the direction opposite to that of the other robot. And, although the self-robot seemed to try to follow the movements of the other-robot, its movements were significantly delayed.

The difference observed between the Wide-IS and Narrow-IS cases is best understood in terms of the different neural dynamic structures developed in these cases. In the case of the probabilistic dynamic structure developed in the Narrow-IS case, the behavior of moving either to the left or to the right is determined simply by following the other-robot by means of sensory reflex without any top-down bias.³⁷ In contrast, in the Wide-IS case, the top-down bias of internal neural activity at decision points is too strong to be modified by sensory input and incorrect movements are initiated and carried through.

3.3 INTRODUCING BOTTOM-UP ERROR REGRESSION

Next, consider an experiment that examines the effects of introducing an additional mechanism of bottom-up error regression into the learned neural dynamics during the course of behavior generation. This is a modified model which maximizes the likelihood LH_{reg} for the time window of the immediate past by modifying the neural activation profile in this past window while fixing the connectivity weights (Fig. 9 (a)) as shown in Eq. (5).

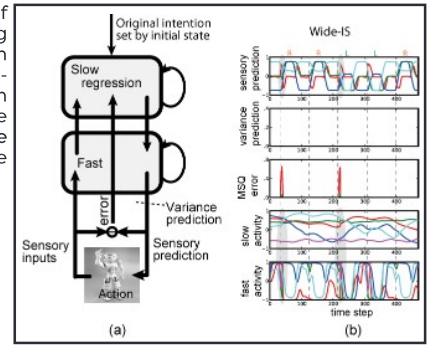
$$\frac{\partial}{\partial \mathbf{w}} \ln L_{reg} = \frac{\partial}{\partial \mathbf{w}} \sum_{t-W}^t \ln \left(\frac{1}{\sigma^2} \exp\left(-\frac{e^2}{2\sigma^2}\right) \right) \quad (5)$$

where the time window is defined from $t-W$ to t at the current time step and the activation states of the slow units at time step $t-W$ (which is the onset of the window) are updated by back propagating the error signal generated. This error regression in terms of updating the activation state at the onset of the window and forward through the window is iterated multiple epochs during each time step in behavior generation. Again, as shown in Eq. 5, error back-propagates more strongly when the estimated variance (as the square error divided by the variance) is smaller.³⁸ An intuitive explanation is that in this scheme the internal representation in the immediate past window is rewritten for the sake of maximizing the likelihood for the ongoing perception.

Fig. 9 (b) shows an example of developments during on-line behavior generation in the trained Wide-IS network using the present error regression scheme.

Clearly, neural activity in the gray area changes in a discontinuous manner with the generation of a sharp peak in prediction error only upon encountering unpredicted action by the other even though this error was rapidly reduced. Note that this sharp peak in prediction error is larger than that generated during on-line prediction in the case of the Narrow-IS as shown in Fig. 8 (b). In short, modulating higher-level neural activity by using error regression caused drastic changes in lower-level network activity including sensory predictions, and in this way prediction errors were rapidly minimized. Ultimately thus, the self-robot was able to re-situate its behavioral context

Figure 9. The results of on-line interaction using the error regression mechanism. (a) The on-line error regression scheme, (b) a sequence generated by the network trained with the Narrow-IS condition.³⁹



immediately after encountering unpredictable events through dense interactions between top-down intentional prediction and bottom-up recognition of actual results.

How can we interpret these experimental results? First, let us summarize what we have just seen. In the Narrow-IS condition, probabilistic network dynamics develop generating actions in a sensory reflex manner. In contrast, proactive behaviors pursuant from deterministic predictions of next actions develop from the Wide-IS condition. It can be said that the Narrow-IS condition develops only weak top-down prior states while the Wide-IS condition develops strong top-down prior states. During the interaction of the self robot with the other robot, the self robot trained under the Narrow-IS condition could easily follow the action sequences arbitrarily determined by the other robot because it simply reacted to sensory inputs, with neural activity at decision points. On the other hand in the Wide-IS condition, the self-robot could not follow the action sequences of the other robot according to sensory inputs, because the top-down bias originating from the initial state was too strong. However, when the error regression scheme was applied utilizing the prediction error generated, the actions of the Wide-IS self robot were modified immediately by means of rapid changes in internal neural states. This bottom-up modulation can be quite strong because the variance is estimated as small in the case of the Wide-IS. This is due to the development of a deterministic dynamic structure, one that plans its next action, and that can be said to “have” a future toward which it has effectively committed itself through proactive cognitive agency given strong top-down prior states. On the other hand, the same force is not so strong in the probabilistic (Narrow-IS) case because the estimated variance at decision points is large. This is to say that the Narrow-IS has plotted no future condition beyond immediate reaction, and has thus cannot be said to “have” a future in this same way.

Consider these Wide and Narrow conditions from the Bayesian viewpoint. In the Bayesian framework, the S-MTRNN represents a likelihood function which maps intention state to a probability distribution of up-coming perceptual states. In these experiments, the distribution of intention states (initial states) was constrained by either the Wide distribution or the Narrow distribution, and the experiments show that the Wide distribution of intention states develops a deterministic dynamics with strong top-down prior states, whereas the Narrow distribution develops a probabilistic process which is a purely reactive process.

3.4 TIME PERCEPTION BY “EMBODIED” RNNs

Now, we come back to the main issue of consciousness. This section briefly looks at the problem of time perception in light of Francisco Varela’s “present-time consciousness.”⁴⁰

Tani and Nolfi postulated that “consciousness” arises at the very moment of segmenting the perceptual flow by means of error regression.⁴¹ Varela’s “present-time consciousness” arises similarly.⁴² First, Varela considered that the immediate past does not belong to a representational conscious memory, but just to an impression consistent with Husserl’s idea of retention.⁴³ So, his question was how the immediate past, experienced just as an impression, could slip into a distant past which can be retrieved through a conscious memory operation later on. And, in response, he proposed that nonlinear dynamics theory could be used as the formal descriptive tool for this phenomenon. By using the phenomenon of the spontaneous flipping of a Necker cube as an example, he explained that the dynamic properties of intermittent chaos characterized by its spontaneous shifts between static and rapid transition modes could explain the paradox of continuous, yet also segmented, time perception.

On his consideration, we may still ask how such spontaneous shifts as those realized by intermittent chaos can be linked to conscious experience. Although Thompson and Varela explain that such shifts are accompanied by shifts in neuronal bias, what is the formal mechanism of this process?⁴⁴ Tani proposes that consciousness arises in the correction and modification of dynamic structures which, in biological cognition, are generated in higher cortical areas.⁴⁵ The following attempts to account for the development of levels of conscious experience in terms of the development of the predictive RNN models described so far in the current paper.

In subjective terms, firstly an agent experiences a continuous perceptual flow without this flow being articulated in any way, that is without this flow representing any discernible thing or event. However, there should be retention and protention in this primordial level, as explained by Husserl (see the last footnote).⁴⁶ Retention and protention are used to designate the experienced sense of the immediate past and the immediate future. They are a part of automatic processes and cannot be controlled consciously. Husserl believed that the subjective experience of “nowness” is extended to include fringes both in the experienced sense of the past and the future in terms of retention and protention. This description of retention and protention at the so-called “pre-empirical” level by Husserl seems to directly correspond to what the basic RNN (as illustrated in Fig. 1 in the earlier section) is performing. The RNN predicts its next state by retaining the past flow in a context dependent way as has been described. This self-organized contextual flow in the forward dynamics of RNNs could account for the phenomenon of retention, whereas prediction based on this contextual flow naturally corresponds to protention.

With Husserl’s idea of “nowness” in terms of retention and protention, the following question arises: Where is the “nowness” bounded? Husserl and Varela believe that the immediate past does not belong to a representational

conscious memory but just to an impression, as suggested above. This led Varela to wonder what kind of mechanism qualitatively changes an experience from just an impression to an episodic consciously retrievable event.⁴⁷ Husserl’s goal was to explain the emergence of objective time from the pre-empirical level of retention and protention,⁴⁸ and he seems to think that the sense of objective time should emerge as a natural consequence of organizing experience into one consistent linear sequence. Still, the question remains: What is the underlying mechanism for this?

One way of approaching this question is to consider first that “nowness” can be bounded where the flow of experience is segmented. Imagine that “Re Fa La” and “Do Mi So” are frequently heard phrases. The sequential notes of “Do Mi So” constitute a chunk within the sound stimulus flow, because the sequence can be predicted perfectly by developing coherence between the predictive neural dynamics and the perceptual flow. Within the chunk, everything proceeds smoothly, automatically, and unconsciously. However, when we hear a next phrase of “Re Fa La” after “Do Mi So” (considering that this second phrase is not necessarily predictable from the first one) a temporal incoherence emerges as prediction error is generated in the transition between the two phrases. The central thesis here is that consciousness arises as the agent attempts to deal with the uncertainty or open possibility between the two.

In Tani and colleagues’ RNN models, the winner module is switched from one to another in MixRNNs or PB is shifted in RNNPB by means of error regression when the external perceptual flow does not match with the internal flow of the prediction. This matching is primarily occurring in the window of the immediate past, as described above. When the prediction is betrayed, the perceptual flow is segmented into chunks associated with shifts of gates or PBs, minimizing prediction error. Those segmented chunks are no longer just parts of the flow, but events that are identified by an activated local module or a PB vector value, e.g., as one of the NAO robot’s behavior primitives. Because of delays in the error minimization process for optimizing gate openings or PB vector, this identification process can be time consuming. This might explain the phenomenological observation that the flow of the immediate past is experienced only as an impression, which later becomes a consciously retrievable object after being segmented. This may correspond to an observation of postdiction evidenced in neuroscience.⁴⁹ See Figure 10 for an illustration of the idea.

The higher level RNN in MixRNNs, RNNPBs, or MTRNNs learns the sequences of the identified events and becomes able to regenerate them as a narrative.

During memory retrieval however, the perceptual flow can be reconstructed only in an indirect way since the flow is now represented by combining a set of commonly used behavior primitives. Although such reconstructions provide for compositionality as well as generalization in representing perceptual flow, they might lose subtle differences or uniqueness in each instance of experience depending on the capacities to retain perceptual dynamics.

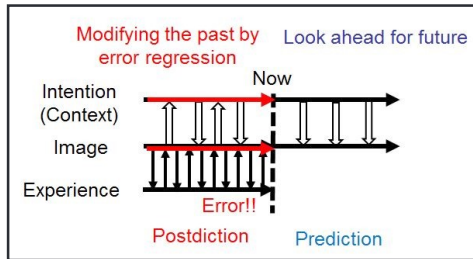


Figure 10. Prediction of future based on the postdiction of the past.

Consequently, we presume that the sense of objective time appears when experience of the perceptual flow is reconstructed as a narrative in a compositional form, while losing its peculiarity.

From the Bayesian perspective of Friston’s FEMP, the agent becomes able to reflect on the episodic sequence with self-estimated certainty when the Narrow-IS condition is applied to S-MTRNN, as shown in the aforementioned experiments by Murata and colleagues.⁵⁰ At this stage, the agent finally becomes able to represent its own episodic sequence in terms of a probabilistic model by inferring that each chunk (moving left or right) simply arises with a certain probability. This is a crucial transition from first reflecting on its own experience as a deterministic “one time only” episodic sequence occurring only in that way, to viewing it as a probability which could have taken place in other ways. From the latter point of view, the agent is successful in ultimately objectifying its own experience by reconstructing it into a generalized model accounting for possible interactions between its self and others. However, it is interesting to note that the agent in this stage does not maintain anymore the subjectivity of naively intending for an uncertain future, because all it maintains is ultimately objectified models of probable futures. Together, these stages of development should begin to account for the process of an agent attaining a reflective self which is only then potentially maintained, for example through inner discourse and conscious narration and which only then results in truly direct subjective experience, the characteristic “mineness” of h-consciousness as revealed in our last paper.

3.5 DISCUSSION

With the composition of intentional sequences, we may understand surprise as their unexpected correction resulting in consciousness. To this, one may object that one becomes conscious of many things without surprise, but this objection is easily answered. Let us consider that intentional processes drive the whole neural network dynamics including the peripheral subnetworks by means of chaos or transient chaos developed in the higher cognitive brain area in order to act on the world in achieving some end of agency such as in Murata’s robot experiment and in Namikawa et al.⁵¹ At this moment of acting, some prediction errors may be generated at the very least because the world is inevitably unpredictable due to its openness relative for instance human cognitive agency. Then, at the very moment when the intention state is modulated by those errors back-propagated from the

peripheral to the higher level, the agent becomes conscious of the formulation of intention upon which it has acted and only in a “postdictive” manner,⁵² i.e., when the intention in the past window is rewritten for the sake of accounting for the current perception, there is consciousness.

With this, we may ask if we can apply the aforementioned analysis to account for the delayed awareness of “free will” as for example evidenced in the famous Libet experiments?⁵³ One might imagine that no prediction errors are to be associated with decisions about pressing a button as in Libet’s experiments. However, in order to initiate a particular movement, internal neural activity in peripheral areas including muscle potential states must prepare for action. With this in mind, prediction errors may arise when the higher cognitive level such as the prefrontal cortex (PFC) or supplementary motor area (SMA) suddenly attempts to drive the lower peripheral processes such as the motor area and somatosensory area through the parietal area, possibly by chaos, to generate a specific movement when the lower parts are not yet prepared for it (see Figure 11).

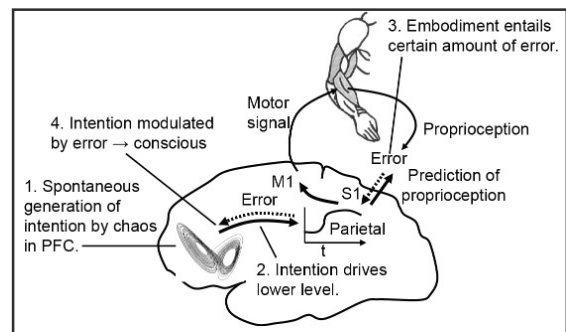


Figure 11. Explanation of how free will can be initiated unconsciously and how one can become consciously aware of it with delay.⁵⁴

In such a situation, a gap may appear between the higher level with the sudden urge for generating the movement and lower level processes which are not yet ready for it. This gap appears in the system as a sort of prediction error, with the intention to act confounded by factors internal to the system as a whole but still external to the intentional processes, themselves. This difference, then, between what is ideally intended and its practical exercise may then cause the conscious awareness of one’s own intention, again with a delay as described by Libet, Gleason and Wright,⁵⁵ and as having been conjectured by Tani.⁵⁶ To sum up, we consider that free will may originate unconsciously by means of the cortical deterministic chaos which can become an object of conscious awareness only after a certain delay under embodiment constraint in terms of postdiction.⁵⁷

Before concluding the current paper, give a close look at the process of error regression at the moment that prediction error increases due to unexpected perception. This error regression process involves the nontrivial phenomena of circular causality, analysis of which reveals subtle characteristics of the conscious process. In simple situations such as shown in the experiments by Murata and

colleagues wherein possible actional decisions are only two, either moving an arm to the left or to the right, the conflictive situation can be resolved instantly by sudden modulation of the intention by error regression.⁵⁸ However, realistic situations are more complex, for example, when a system has to perform online modification of a goal-directed plan by searching among various possible combinations of behavior primitives by means of error regression, while retaining immediate integrity in the face of environmental forces, i.e., adapting to rapid changes of the current situation until the newly formulated intention can be carried forward.

4. CONCLUSION

The current paper reviewed a series of neurobotics studies conducted by Jun Tani and colleagues that attempt to provide a purely formal, structural account of dynamical processes essential for consciousness. The core ingredients of Tani's models are prediction and postdiction through predictive coding and implemented in different recurrent neural network (RNN) models that together represent a progression from reflexive to proactive, self-reflective and creative agency. The review moved from simple to more complex model hierarchies.

Robotics experiments employing these models clarified dynamics inherent in levels of consciousness from momentary self-consciousness (surprise) to narrative self and reflective self-consciousness (the "chunking" of experience and the articulation of perceptual flow according to developing action potentials). The paper concluded with a brief phenomenological analysis of time perception within this family of models, including model extensions accounting for free will and its characteristic postdictive conscious awareness. In the next paper, we will begin with some of Tani and colleagues' work on these model extensions into more complex situations, before returning to Boltuc's naturalistic non-reductionism and a philosophical analysis of any claim to consciousness of artificial systems.

NOTES

1. The correspondence should be sent to tani1216jp@gmail.com
2. Here it is interesting to note that predictive coding is inspired by studies on biological systems, so computational architectures employing predictive coding are by definition instances of a biological approach. R. N. Rao and D. H. Ballard, "Predictive Coding in the Visual Cortex: A Functional Interpretation of Some Extra-Classical Receptive-Field Effects," *Nature Neuroscience* 2 (1999): 79–87; J. Tani and S. Nolfi, "Learning to Perceive the World as Articulated: An Approach for Hierarchical Learning in Sensory-Motor Systems," *Neural Networks* 12, no. 7 (1999): 1131–41; K. Friston, "A Theory of Cortical Responses," *Philosophical Transactions of the Royal Society B: Biological Sciences* 360, no. 1456 (2005): 815–36.
3. J. Tani, "Autonomy of 'Self' at Criticality: The Perspective from Synthetic Neuro-Robotics," *Adaptive Behavior* 17, no. 5 (2009): 421–43; A. Clark, *Surfing Uncertainty: Prediction, Action, and the Embodied Mind* (NY: Oxford University Press, 2015); J. Tani, *Exploring Robotic Minds: Actions, Symbols, and Consciousness as Self-Organizing Dynamic Phenomena* (New York: Oxford University Press, 2016).
4. J. Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process," *Neural Networks* 16 (2003): 11–23.

5. Friston, "A Theory of Cortical Responses"; K. Friston, "The Free-Energy Principle: A Unified Brain Theory?" *Nature Reviews Neuroscience* 11 (2010): 127–38.
6. M. I. Jordan, "Serial Order: A Parallel Distributed Processing Approach," Technical Report, California University, San Diego, 1986; J. L. Elman, "Finding Structure in Time," *Cognitive Science* 14 (1990): 179–211; R. J. Williams and D. Zipser, "Finding Structure in Time," Institute for Cognitive Science Report, University of California, San Diego, 1990.
7. Ibid.
8. C.f. J. Tani and N. Fukumura, "Embedding a Grammatical Description in Deterministic Chaos: An Experiment in Recurrent Neural Learning," *Biological Cybernetics* 72, no. 4 (1995): 365–70.
9. D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning Internal Representations By Error Propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, ed. D. E. Rumelhart and J. L. McClelland (Cambridge, MA: The MIT Press, 1986).
10. J. Tani and S. Nolfi, "Self-Organization of Modules and Their Hierarchy in Robot Learning Problems: A Dynamical Systems Approach," *News Letter on System Analysis for Higher Brain Function Research Project* 2, no. 4 (1997): 1–11; Tani and Nolfi, "Learning to Perceive the World as Articulated: An Approach for Hierarchical Learning in Sensory-Motor Systems."
11. Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process"; J. Tani and M. Ito, "Self-Organization of Behavioral Primitives as Multiple Attractor Dynamics: A Robot Experiment," *Systems, Man, and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, 33, no. 4 (2003): 481–88.
12. Rizzolatti et al. ("Mirror Neuron: A Neurological Approach to Empathy," in *Neurobiology of Human Values* [Springer-Verlag Berlin Heidelberg, 1995]) as explored in J. Tani, M. Ito, and Y. Sugita, "Self-Organization of Distributedly Represented Multiple Behavior Schemata in a Mirror System: Reviews of Robot Experiments Using RNNPB," *Neural Networks* 17 (2004): 1273–89.
13. Tani and Nolfi, "Self-Organization of Modules and Their Hierarchy in Robot Learning Problems"; Tani and Nolfi, "Learning to Perceive the World as Articulated."
14. Here, recall that it is the object of the network to minimize error through the error back-propagation through time (BPTT) algorithm.
15. Tani and Nolfi, "Learning to Perceive the World as Articulated."
16. Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process"; Tani and Ito, "Self-Organization of Behavioral Primitives as Multiple Attractor Dynamics: A Robot Experiment."
17. Tani and Nolfi, "Self-Organization of Modules and Their Hierarchy in Robot Learning Problems"; Tani and Nolfi, "Learning to Perceive the World as Articulated"; and Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process."
18. It is interesting to note here a formal similarity with "abductive" agent-level evolutionary models.
19. Tani and Nolfi, "Self-Organization of Modules and Their Hierarchy in Robot Learning Problems"; Tani and Nolfi, "Learning to Perceive the World as Articulated."
20. Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process."
21. Y. Yamashita and J. Tani, "Emergence of Functional Hierarchy in a Multiple Timescale Neural Network Model: A Humanoid Robot Experiment," *PLoS Computational Biology* 4, no. 11 (2008): e1000220.
22. Ibid.
23. R. Nishimoto and J. Tani, "Development of Hierarchical Structures for Actions and Motor Imagery: A Constructivist View from Synthetic Neuro-Robotics Study," *Psychological Research* 73, no. 4 (2009): 545–58; and H. Arie, T. Endo, T. Arakaki, S. Sugano, and J. Tani, "Creating Novel Goal-Directed Actions at Criticality: A Neuro-robotic Experiment," *New Mathematics and Natural Computation* 5, no. 01 (2009): 307–34.

24. D. T. Campbell, "'Downward Causation' in Hierarchically Organized Biological Systems," in *Studies in the Philosophy of Biology* (Macmillan Education UK, 1974), 179–86; E. Thompson and F. J. Varela, "Radical Embodiment: Neural Dynamics and Consciousness," *Trends in Cognitive Sciences* 5, no. 10 (2001): 418–25.

25. Friston, "A Theory of Cortical Responses"; K. Friston, "Hierarchical Models in the Brain," *PLoS Computational Biology* 4, no. 11 (2008): e1000211.

26. It is important here to bear in mind that these are systems enabling agency, and so an action that ends very far from a target is much worse than one which ends close enough. It is not as innocuous as simply getting something wrong. If variance is high, then a prediction which hits its target is extremely accurate, such that in the real world it may not be believed, e.g., "too good to be true." However, when variance is high, it also means that average values do not effectively inform action. Acting on the basis of an average will always in the long run result in error proportional to variance. Once this is understood, then an agent may apply estimated variance in the prediction of optimal next actions, as this value may inform the agent what to expect given prior instances, reducing error over the long run.

This formal model recalls Plato's concern with the science of science that is ultimately knowledge of good and bad, a second-order understanding that for example directs sight but never sees a thing, c.f. Charmides.

27. Friston, "A Theory of Cortical Responses."

28. Friston, "The Free-Energy Principle: A Unified Brain Theory?"

29. S. Murata, Y. Yamashita, H. Arie, T. Ogata, S. Sugano, and J. Tani, "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others: A Neuro-robotics Experiment," *IEEE Trans. on Neural Networks and Learning Systems*. 2015. doi:10.1109/TNNLS.2015.2492140

30. Recognition density in Friston, "The Free-Energy Principle: A Unified Brain Theory?"

31. Redrawn from Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

32. Ibid.

33. C.f. Yamashita and Tani, "Emergence of Functional Hierarchy in a Multiple Timescale Neural Network Model: A Humanoid Robot Experiment"; J. Namikawa, R. Nishimoto, and J. Tani, "A Neurodynamic Account of Spontaneous Behavior," *PLoS Computational Biology* 7, no. 10 (2011): e1002221.

34. Redrawn from Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

35. Namikawa et al., "A Neurodynamic Account of Spontaneous Behavior."

36. Redrawn from Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

37. Recall that the slow and fast networks showed the same dynamics at each point.

38. In terms of human experience, it feels worse being wrong when sure that he/she is right than when it is a recognized matter of chance.

39. Redrawn from Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

40. F. J. Varela, "Present-Time Consciousness," *Journal of Consciousness Studies* 6, nos. 2-3 (1999): 111–40.

41. Tani and Nolfi, "Learning to Perceive the World as Articulated."

42. Varela, "Present-Time Consciousness."

43. E. Husserl, "The Phenomenology of Internal Time Consciousness," trans. J. S. Churchill (Bloomington, IN: Indiana University Press, 1964). Husserl introduced the famous idea of "retention" and "protention" for explaining the paradoxical nature of "nowness." He used an example of hearing a sound phrase such as "Do Mi So" for explaining the idea. When we hear the note "Mi," we would still perceive a lingering impression of "Do," and at the same time we would anticipate hearing the next note of "So." The former is called retention and the latter protention. These terms

are used to designate the experienced sense of the immediate past and the immediate future.

44. E. Thompson and F. J. Varela, "Radical Embodiment: Neural Dynamics and Consciousness," *Trends in Cognitive Sciences* 5, no. 10 (2001): 418–25.

45. J. Tani, "An Interpretation of the 'Self' from the Dynamical Systems Perspective: A Constructivist Approach," *Journal of Consciousness Studies* 5, nos. 5/6 (1998): 516–42; Tani and Nolfi, "Learning to Perceive the World as Articulated: An Approach for Hierarchical Learning in Sensory-Motor Systems"; Tani, "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process."

46. Husserl, "The Phenomenology of Internal Time Consciousness." See footnote 51.

47. Varela, "Present-Time Consciousness."

48. Husserl, "The Phenomenology of Internal Time Consciousness."

49. D. M. Eagleman and T. J. Sejnowski, "Motion Integration and Postdiction in Visual Awareness," *Science* 287, no. 5460 (2000): 2036–38; S. Shimojo, "Postdiction: Its Implications on Visual Awareness, Hindsight, and Sense of Agency," *Frontiers in Psychology* 5 (2014): 196.

50. Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

51. Ibid. and in Namikawa et al., "A Neurodynamic Account of Spontaneous Behavior."

52. Shimojo, "Postdiction: Its Implications on Visual Awareness, Hindsight, and Sense of Agency."

53. C.f. B. Libet, E. W. Wright, and C. A. Gleason, "Preparation- or Intention-to-Act, in Relation to Pre-event Potentials Recorded at the Vertex," *Electroencephalography and Clinical Neurophysiology* 56, no. 4 (1983): 367–72; B. Libet, "Unconscious Cerebral Initiative and the Role of Conscious Will in Voluntary Action," *Behavioral and Brain Sciences* 8 (1985): 529–39.

54. Redrawn from J. Tani, *Exploring Robotic Minds: Actions, Symbols, and Consciousness as Self-Organizing Dynamic Phenomena* (New York: Oxford University Press, 2016) with permission from Oxford University Press.

55. Libet, Gleason, and Wright, "Preparation- or Intention-to-Act, in Relation to Pre-event Potentials Recorded at the Vertex." The preceding interpretation accords with recent work in the sense of self-agency as it changes due to feedback. In short, the more that an agent perceives itself to be in control of outcomes, the more it feels the sense of ownership of actions performed (c.f. N. Kumar, J. A. Manjaly, and K. P. Miyapuram, "Feedback about Action Performed Can Alter the Sense of Self-Agency," *Frontiers in Psychology*, February 25, 2014).

56. Tani, *Exploring Robotic Minds: Actions, Symbols, and Consciousness as Self-Organizing Dynamic Phenomena*.

57. Eagleman and Sejnowski, "Motion Integration and Postdiction in Visual Awareness"; Shimojo, "Postdiction: Its Implications on Visual Awareness, Hindsight, and Sense of Agency."

58. Murata et al., "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others."

REFERENCES

Arie, H., T. Endo, T. Arakaki, S. Sugano, and J. Tani. "Creating Novel Goal-Directed Actions at Criticality: A Neuro-robotic Experiment." *New Mathematics and Natural Computation* 5, no. 01 (2009): 307–34.

Boltuc, P. "The Philosophical Issue in Machine Consciousness." *International Journal of Machine Consciousness* 1, no. 1 (2009): 155–76.

Campbell, D. T. "'Downward Causation' in Hierarchically Organized Biological Systems." In *Studies in the Philosophy of Biology*, 179–86. Macmillan Education UK, 1974.

Clark, A. *Surfing Uncertainty: Prediction, Action, and the Embodied Mind*. NY: Oxford University Press, 2015.

Eagleman, D. M., and T. J. Sejnowski. "Motion Integration and Postdiction in Visual Awareness." *Science* 287, no. 5460 (2000): 2036–38.

Elman, J. L. "Finding Structure in Time." *Cognitive Science* 14 (1990): 179–211.

- Friston, K. "A Theory of Cortical Responses." *Philosophical Transactions of the Royal Society B: Biological Sciences* 360, no. 1456 (2005): 815–36.
- . "Hierarchical Models in the Brain." *PLoS Computational Biology* 4, no. 11 (2008): e1000211.
- . "The Free-Energy Principle: A Unified Brain Theory?" *Nature Reviews Neuroscience* 11 (2010): 127–38.
- Husserl, E. "The Phenomenology of Internal Time Consciousness." Translated by J. S. Churchill. Bloomington, IN: Indiana University Press, 1964.
- Jordan, M. I. "Serial Order: A Parallel Distributed Processing Approach." Technical Report. California University, San Diego, 1986.
- Kumar, N., J. A. Manjaly, and K. P. Miyapuram. "Feedback about Action Performed Can Alter the Sense of Self-Agency." *Frontiers in Psychology*, February 25, 2014. doi:10.3389/fpsyg.2014.00145
- Libet, B. "Unconscious Cerebral Initiative and the Role of Conscious Will in Voluntary Action." *Behavioral and Brain Sciences* 8 (1985): 529–39.
- Libet, B., E. W. Wright, and C. A. Gleason. "Preparation- or Intention-to-Act, in Relation to Pre-event Potentials Recorded at the Vertex." *Electroencephalography and Clinical Neurophysiology* 56, no. 4 (1983): 367–72.
- Murata, S., Y. Yamashita, H. Arie, T. Ogata, S. Sugano, and J. Tani. "Learning to Perceive the World as Probabilistic or Deterministic via Interaction with Others: A Neuro-robotics Experiment." *IEEE Trans. on Neural Networks and Learning Systems*. 2015. doi:10.1109/TNNLS.2015.2492140
- Namikawa, J., R. Nishimoto, and J. Tani. "A Neurodynamic Account of Spontaneous Behavior." *PLoS Computational Biology* 7, no. 10 (2011): e1002221.
- Nishimoto, R., and J. Tani. "Development of Hierarchical Structures for Actions and Motor Imagery: A Constructivist View from Synthetic Neuro-Robotics Study." *Psychological Research* 73, no. 4 (2009): 545–58.
- Rao, R. N., and D. H. Ballard. "Predictive Coding in the Visual Cortex: A Functional Interpretation of Some Extra-Classical Receptive-Field Effects." *Nature Neuroscience* 2 (1999): 79–87.
- Reed, E. S., and D. Schoenherr. "The Neuropathology of Everyday Life: On the Nature and Significance of Microslips in Everyday Activities." Unpublished manuscript. 1992.
- Rumelhart, D. E., G. E. Hinton, and R. J. Williams. "Learning Internal Representations By Error Propagation." In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, edited by D. E. Rumelhart and J. L. McClelland. Cambridge, MA: The MIT Press, 1986.
- Shimojo, S. "Postdiction: Its Implications on Visual Awareness, Hindsight, and Sense of Agency." *Frontiers in Psychology* 5 (2014): 196. 10.3389/fpsyg.2014.00196.
- Tani, J. "An Interpretation of the 'Self' from the Dynamical Systems Perspective: A Constructivist Approach." *Journal of Consciousness Studies* 5, nos. 5/6 (1998): 516–42.
- . "Learning to Generate Articulated Behavior Through the Bottom-Up and the Top-Down Interaction Process." *Neural Networks* 16 (2003): 11–23.
- . "The Dynamical Systems Accounts for Phenomenology of Immanent Time: An Interpretation By Revisiting a Robotics Synthetic Study." *Journal of Consciousness Studies* 11, no. 9 (2004): 5–24.
- . "Autonomy of 'Self' at Criticality: The Perspective from Synthetic Neuro-Robotics." *Adaptive Behavior* 17, no. 5 (2009): 421–43.
- . *Exploring Robotic Minds: Actions, Symbols, and Consciousness as Self-Organizing Dynamic Phenomena*. New York: Oxford University Press, 2016.
- Tani, J., and N. Fukumura. "Embedding a Grammatical Description in Deterministic Chaos: An Experiment in Recurrent Neural Learning." *Biological Cybernetics* 72, no. 4 (1995): 365–70.
- Tani, J., and S. Nolfi. "Self-Organization of Modules and Their Hierarchy in Robot Learning Problems: A Dynamical Systems Approach." *News Letter on System Analysis for Higher Brain Function Research Project 2*, no. 4 (1997): 1–11.
- . "Learning to Perceive the World as Articulated: An Approach for Hierarchical Learning in Sensory-Motor Systems." *Neural Networks* 12, no. 7 (1999): 1131–41.
- Tani, J., and M. Ito. "Self-Organization of Behavioral Primitives as Multiple Attractor Dynamics: A Robot Experiment." *Systems, Man, and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, 33, no. 4 (2003): 481–88.
- Tani, J., M. Ito, and Y. Sugita. "Self-Organization of Distributedly Represented Multiple Behavior Schemata in a Mirror System: Reviews of Robot Experiments Using RNNPB." *Neural Networks* 17 (2004): 1273–89.
- Thompson, E., and F. J. Varela. "Radical Embodiment: Neural Dynamics and Consciousness." *Trends in Cognitive Sciences* 5, no. 10 (2001): 418–25.
- Varela, F. J. "Present-Time Consciousness." *Journal of Consciousness Studies* 6, nos. 2-3 (1999): 111–40.
- Williams, R. J., and D. Zipser. "Finding Structure in Time." Institute for Cognitive Science Report. University of California, San Diego, 1990.
- Yamamoto, J., J. Suh, D. Takeuchi, and S. Tonegawa. "Successful Execution of Working Memory Linked to Synchronized High-Frequency Gamma Oscillations." *Cell* 157, no. 4 (2014): 845–57.
- Yamashita, Y., and J. Tani. "Emergence of Functional Hierarchy in a Multiple Timescale Neural Network Model: A Humanoid Robot Experiment." *PLoS Computational Biology* 4, no. 11 (2008): e1000220.

Kant on Constituted Mental Activity

Richard Evans
IMPERIAL COLLEGE, UK

1 INTRODUCTION

Consider the following functionalist claim:

There is an architecture, describable in the language of computer science, such that any creature or machine that realises this architecture thereby counts as a cognitive agent, an agent with original (non-derivative) intentionality.

Some of the more practically minded among us will be dissatisfied with this existentially quantified assertion: rather than just saying that there is *some* such architecture, it would be much more helpful to know exactly what this architecture is. What sort of architecture could satisfy such a claim?

I believe the answer to this question has been hiding in plain sight for over two hundred years: in *The Critique of Pure Reason*, Kant provides a detailed description of just such an architecture.

At the heart of Kant's vision is the *self-legislating agent*: an agent who constructs rules that he then solemnly follows. The Kantian cognitive architecture is a particular type of *computational process*: a rule-induction process. If this rule-induction process satisfies certain constraints, then—Kant claims—the process' internal activities count as cognitive activities.

This paper sketches the philosophical background behind this architecture. It attempts to motivate and defend Kant's vision of a self-legislating computational agent.¹

2 MENTAL ACTIVITY AS CONSTITUTED ACTIVITY

We are familiar with the idea that social activity is *constituted* activity. Pushing the wooden horse-shaped piece forward counts, in the right circumstances, as moving the knight to