

Investigation of the Sense of Agency in Social Cognition, based on frameworks of Predictive Coding and Active Inference: A simulation study on multimodal imitative interaction

Wataru Ohata and Jun Tani *

Cognitive Neurorobotics Research Unit, Okinawa Institute of Science and Technology Graduate University, Okinawa, Japan

Correspondence*:
Jun Tani
jun.tani@oist.jp

2 ABSTRACT

3 When agents interact socially with different intentions (or wills), conflicts are difficult to avoid.
4 Although the means by which social agents can resolve such problems autonomously has not
5 been determined, dynamic characteristics of agency may shed light on underlying mechanisms.
6 Therefore, the current study focused on the sense of agency, a specific aspect of agency referring
7 to congruence between the agent's intention in acting and the outcome, especially in social
8 interaction contexts. Employing predictive coding and active inference as theoretical frameworks
9 of perception and action generation, we hypothesize that regulation of complexity in the evidence
10 lower bound of an agent's model should affect the strength of the agent's sense of agency and
11 should have a significant impact on social interactions. To evaluate this hypothesis, we built a
12 computational model of imitative interaction between a robot and a human via visuo-proprioceptive
13 sensation with a variational Bayes recurrent neural network, and simulated the model in the form
14 of pseudo-imitative interaction using recorded human body movement data, which serve as the
15 counterpart in the interactions. A key feature of the model is that the complexity of each modality
16 can be regulated differently by changing the values of a hyperparameter assigned to each local
17 module of the model. We first searched for an optimal setting of hyperparameters that endow
18 the model with appropriate coordination of multimodal sensation. These searches revealed that
19 complexity of the vision module should be more tightly regulated than that of the proprioception
20 module because of greater uncertainty in visual information flow. Using this optimally trained
21 model as a default model, we investigated how changing the tightness of complexity regulation
22 in the entire network after training affects the strength of the sense of agency during imitative
23 interactions. The results showed that with looser regulation of complexity, an agent tends to
24 act more egocentrically, without adapting to the other. In contrast, with tighter regulation, the
25 agent tends to follow the other by adjusting its intention. We conclude that the tightness of
26 complexity regulation significantly affects the strength of the sense of agency and the dynamics
27 of interactions between agents in social settings.

28 **Keywords:** sense of agency, predictive coding, active inference, multimodal perception, human-robot interaction, recurrent neural
29 network, variational Bayes

1 INTRODUCTION

30 Humans are social beings by nature, and each individual regularly interacts with others in various ways.
31 Even though individuals act based on their intentions or wills, they sometimes acts in agreement with
32 others, doing something collaboratively, while at other times they disagree. Either case may be conscious or
33 unconscious. What determines such the type of interaction and how? To evaluate this problem, we consider
34 possible relationships between *agency* of each individual and social interactions between individuals.
35 Then, we introduce predictive coding and active inference to formulate the problem in a computational
36 framework and we propose a specific hypothesis to predict the type of interaction. We deliver a schematic
37 of our computational model and experimental setup to evaluate the hypothesis and conclude the section by
38 highlighting some critical findings.

39 1.1 Agency in social cognition

40 In social interactions, agents sometimes cooperate by sharing intentions so as to derive mutual benefits,
41 while at other times they cause conflicts by following their own intentions and ignoring the interests of
42 others. Although how such complexities in social interactions emerge is not obvious, we hypothesized
43 that dynamic characteristics of *agency* in social interactions might shed light on underlying mechanisms.
44 Recently, the study of agency has attracted considerable attention from researchers in various disciplines,
45 including philosophy, psychology, cognitive science, and neuroscience. Specifically, the sense of agency
46 (SoA) (Gallagher, 2000; Synofzik et al., 2008; Moore et al., 2009) refers to congruence between an
47 agent's intention or belief in an action and its anticipated outcome, which endows the agent with the
48 sense that "*I am the one generating this action*". Along with studies in experimental psychology, building
49 a computational model of SoA is also important in order to explore the nature of agency (Legaspi and
50 Toyozumi, 2019). In the study of computational models of agents, predictive coding (PC) (Rao and Ballard,
51 1999; Tani and Nolfi, 1999; Lee and Mumford, 2003; Friston, 2005; Hohwy, 2013; Clark, 2015; Friston,
52 2018) and active inference (AIF) (Friston et al., 2009, 2010; Baltieri and Buckley, 2017; Buckley et al.,
53 2017; Pezzulo et al., 2018; Oliver et al., 2019) have recently attracted considerable attention since they
54 provide rigid theoretical frameworks for defining perception and action generation. In the framework of
55 PC and AIF, an agent's intention or belief can be formulated as a predictive model, and it is thought that
56 congruence between the prediction of action outcomes and observations reinforces the SoA (Friston, 2012).

57 In situations involving social interaction, however, where multiple agents interact, it becomes challenging
58 for each agent to sustain its SoA, because other agents, having their own intentions, may not act as desired.
59 If social agents are required to coordinate actions so as to obtain benefits by minimizing possible conflicts,
60 we speculated that the strength of agency should be arbitrated among those agents during some conflicts.
61 Let us consider a dyadic synchronized imitation as an example of social interaction, wherein two agents
62 attempt to synchronously imitate one another's movement patterns using predictions based on prior learning.
63 In addition, let us assume a setting in which two agents imitate one another in sequences of movement
64 patterns based on memorized transition rules, in which unpredictable transitions in movement patterns are
65 included. For example, either movement pattern B or C can appear after movement pattern A (see also
66 Figure 3.2 (A)). In this setting, agent 1 may opt for movement pattern B after A, acting as a leader with
67 *strong agency* and agent 2 may just follow agent 1 by generating pattern B with *weak agency*. This can
68 result in successful mutual imitation without generating conflict. However, if both agents maintain strong
69 agency, each may generate its own pattern (B or C) without compromise, resulting in conflict.

70 While investigating agency in social interactions, we concluded that it would also be worthwhile to
71 consider how agency and mirror neuron systems (MNS) (Rizzolatti and Fogassi, 2014; Kilner et al.,
72 2007) might be related, since MNS are thought to contribute to various types of social cognitive behavior,

73 including imitation (Hurley, 2005). MNS was first discovered in area F5 of the monkey premotor cortex
 74 (Di Pellegrino et al., 1992; Gallese et al., 1996), and it is activated when monkeys execute their own actions,
 75 as well as when observing those performed by others. Because MNS uses observations of an action to
 76 generate the same action, it may participate in imitative behaviors, which are thought to be the basis of
 77 various higher cognitive functions (Aly and Tapus, 2015; Kohler et al., 2002; Oztop et al., 2006, 2013). A
 78 natural question regarding such an MNS mechanism is how the agency of each individual can be exerted
 79 if MNS is the default mode. Intention to generate an action could conflict with an automatic response to
 80 imitate an action demonstrated by others. Although modeling studies of MNS have also been conducted
 81 from the view point of PC and AIF using Bayesian frameworks (Friston et al., 2011; Kilner et al., 2007)
 82 and by using deterministic recurrent neural networks (RNNs) (Ito and Tani, 2004; Ahmadi and Tani, 2017;
 83 Hwang et al., 2020), the aforementioned problem of agency has not been well considered.

84 1.2 Predictive coding and active inference

85 Next, let us consider how the strength of agency can be modeled using a framework of PC and AIF. For
 86 this purpose, first we briefly review the concepts of PC and their mathematical properties, as follows. In PC,
 87 perception is thought to be achieved via iterative interactions between a prior expectation of a sensation
 88 and a posterior inference from a sensory outcome. The prior expectation of the sensation can be modeled
 89 by statistical generative models that map the prior of the latent variable to the sensory expectation. The
 90 posterior inference can be carried out by taking the error between the expected sensation and its outcome
 91 and by updating the posterior of the latent variable in the direction of minimizing the error, under the
 92 constraint of minimizing Kullback-Leibler divergence (KL divergence) between the posterior distribution
 93 and that of the prior. Typically, both the prior and the posterior are represented by Gaussian distributions
 94 with parameters of mean and variance, as will be described later. This is equal to maximizing the lower
 95 bound of the logarithm of marginal likelihood (a.k.a evidence lower bound) expressed by two terms:
 96 *accuracy* and *complexity*.

$$97 \quad \ln p_{\theta}(\mathbf{X}) \geq \underbrace{\int q_{\phi}(\mathbf{z}|\mathbf{X}) \ln \frac{p_{\theta}(\mathbf{X}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{X})} dz}_{\text{Evidence lower bound}} \quad (1)$$

$$98 \quad = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{X})}[\ln p_{\theta}(\mathbf{X}|\mathbf{z})]}_{\text{Accuracy}} - \underbrace{D_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{X})||p(\mathbf{z})]}_{\text{Complexity}} \quad (2)$$

99 where \mathbf{X} is the observation, \mathbf{z} is the latent variable, $q_{\phi}(\mathbf{z}|\mathbf{X})$ is the approximate posterior, and θ and ϕ are
 100 the parameters of the model. *Accuracy* is the expectation of log-likelihood with respect to the approximate
 101 posterior, which represents reconstruction of the observation with the approximate posterior. *Complexity* is
 102 the KL divergence between the approximate posterior and the prior, which serves to regularize the model.
 103 Importantly, in maximizing the lower bound, the interplay between these two terms characterizes how
 104 the model behaves in learning and prediction (Higgins et al., 2017). Maximization of the lower bound is
 105 equivalent to minimization of *free-energy* proposed by Friston (Friston, 2005).

106 Next, AIF is described briefly. AIF explains that action or motor commands should be generated so that
 107 their sensory outcomes coincide with expected outcomes. As a simple example, consider how expected
 108 proprioception in terms of robot joint angles can be achieved by generating sufficient motor torque. This
 109 can be done with an inverse model that maps expected joint angles to the required motor torques, or by
 110 employing a PID controller such that necessary motor torque to minimize errors between expected joint
 111 angles and actual angles can be derived by means of a simple error feedback mechanism. Both PC and AIF
 112 attempt to minimize error between the expected sensation and the actual outcome; however, in PC this is
 113 accomplished by changing the intention via the posterior inference and by changing the environment state

114 through action in AIF. When PC and AIF are performed in tandem, while an agent acts on the environment,
115 an agent with a more precise prior (smaller variance) should behave with strong agency, being less likely to
116 change its own intention, and more likely to change the environmental state. On the other hand, an agent
117 with a less precise prior (with larger variance) should behave with weaker agency, being more likely to
118 change its own intention than the environmental state.

119 **1.3 Related work**

120 Although PC and AIF have attracted much attention from brain modeling researchers, it is unusual to
121 see them used in computational studies employing learnable neural network models, especially those
122 that can handle continuous spatio-temporal patterns characterized in multimodal sensory inputs. To this
123 intent, Ahmadi and Tani (Ahmadi and Tani, 2019) recently proposed so-called, Predictive-coding-inspired
124 Variational Recurrent Neural Network (PV-RNN). PV-RNN is a type of variational recurrent neural network
125 that approximates the posterior at each time step in sequential patterns with variational inference, and is
126 formalized by employing predictive coding. By making predictions in the form of the sequential prior
127 (Chung et al., 2015) with time-varying parameterized Gaussian distribution, PV-RNN enables the model
128 to represent strength of intention or agency. Ahmadi and Tani (2019) introduced a hyper parameter w
129 called the meta-prior, which weights regulation of the complexity term in the evidence lower bound (the
130 second term in equation 2). They found that a model trained with looser regulation of the complexity term,
131 achieved by setting the meta-prior to a larger value, develops more deterministic dynamics with higher
132 estimated precision in the sequential prior, whereas a model trained with tighter regulation, accomplished
133 by setting the meta-prior to a smaller value, develops more probabilistic dynamics with lower estimated
134 precision. In another attempt to implement free-energy minimization with an artificial neural network, Pitti
135 et al. (Pitti et al., 2020) proposed a spiking neural network architecture that minimizes free-energy to model
136 the fronto-striatal system in the brain.

137 Chame and Tani (Chame and Tani, 2019) used PV-RNN to conduct a human-robot interaction experiment
138 using a single perceptual channel of proprioception. Although their analysis of the experiments was
139 preliminary, they suggested that when the model is trained under looser regulation of the complexity
140 term, the model tends to behave egocentrically, adapting less to proprioceptive inputs, whereas under
141 tighter regulation of the complexity term, the network tends to behave more passively, adapting more
142 to proprioceptive inputs. However, such network characteristics, once developed through learning under
143 particular conditions to regulate the complexity term, cannot be changed thereafter. In social interactions, it
144 is natural that agents act differently, depending on the social context at a given moment. Sometimes they
145 tend to preserve their prior intention by acting perversely, and at other times they change it more easily by
146 adapting to intentions of others. The current study examines whether such shifts in strength of agency can
147 be achieved during the interaction phase by changing the value of the meta-prior from the default strength
148 set in the learning phase.

149 **1.4 Imitative interaction using a variational Bayes recurrent neural network**

150 Here, we explain the general concept underlying our computational model, experimental design, and
151 obtained results. We first proposed an artificial neural network model that can be applied to an imitative
152 interaction task using multimodal sensation of vision and proprioception by extending PV-RNN. PV-RNN
153 is used because to our knowledge this network model is the only RNN-type model that can instantiate
154 predictive coding and active inference in a continuous spatio-temporal domain by following a Bayesian
155 framework. The proposed model is comprised of a multi-layered PV-RNN with a branching structure, in
156 which two branches responsible for perception of vision and proprioception are connected through an
157 associative module. In addition, the current model inherits the structure of Multiple Timescale Recurrent

158 Neural Network (MTRNN) (Yamashita and Tani, 2008). MTRNN extracts a temporal hierarchy contained
159 in sequential patterns (Yamashita and Tani, 2008; Nishimoto and Tani, 2009; Hwang et al., 2020). By
160 assigning faster timescales to the peripheral sensory modules for vision and proprioception and slower
161 timescales to the associative module, hierarchical multimodal integration from sensory-motor levels to
162 abstract intention levels should be achieved.

163 The entire network model is considered a generative model that predicts incoming visual sensation and
164 proprioception simultaneously through a generative process along with a top-down pathway from the
165 associative module to both of the sensory modules. The resultant prediction error for each sensory modality
166 is back-propagated through time (Werbos, 1974; Rumelhart et al., 1985) (BPTT) and through each module
167 to the associative module, by which the latent state in each module is modulated so as to maximize the
168 evidence lower bound shown in the equation 2. This corresponds to the posterior inference. The network is
169 trained through supervised learning by maximizing the evidence lower bound.

170 However, coordinating multimodal sensations appropriately is still not an easy problem when intrinsic
171 complexity and randomness in spatio-temporal patterns differ in each modality (Ogata et al., 2010; Valentin
172 et al., 2019). Studies on cue integration in multimodal sensation have shown that inferences about the
173 hidden state of the environment should be accomplished by assigning the greatest weight to information
174 obtained from the most reliable sensory modality (Battaglia et al., 2003). In predictive coding, reliability
175 can be represented by the accuracy estimated for each modality of the sensory model, provided that its
176 generalization is preserved by minimizing model complexity adequately when the amount of training data
177 is limited. We speculate that the complexity term should be regulated adequately for each sensory modality
178 during training, such that the best generalization can be achieved for each. Since each PV-RNN module can
179 be assigned different values of the meta-prior, the above could be achieved by searching for an adequate
180 value of each meta-prior through trial and error during the learning phase.

181 The proposed model was evaluated by simulating “pseudo” imitative interaction using visuo-
182 proprioceptive sequence patterns recorded from human demonstrators. Although human-robot interaction
183 should be studied in a physical system to allow the human and the robot to respond to each other in an
184 online fashion, it is difficult for the current system to work in real time because inference of the posterior
185 using PV-RNN is computationally intensive, especially when pixel-level vision is used as one of the sensory
186 modalities. Therefore, the current study focuses on simulation experiments using pre-recorded data.

187 First, we investigated how changing the tightness used to regulate the complexity term for each sensory
188 module in the learning phase affects the quality of integrating multimodal sensation in an imitative
189 interaction. For this purpose, we examined possible effects of assigning different values of the meta-prior to
190 the vision module and the proprioceptive module, on performance characteristics in learning, as well as in
191 the resulting imitative interaction. Our results suggest that regulating complexity more in the vision module
192 than in the proprioception module facilitates better imitation performance in multi-modal sensation after
193 learning, because visual sensory information contains more randomness than proprioceptive information.

194 Second, as the main motivation of the current study, we investigated how changing the tightness used
195 to regulate the complexity term in the entire network after the learning phase affects the strength of
196 agency. Using a network trained by tuning the meta-priors assigned to each sensory module in the previous
197 experiment, we examined how increasing or decreasing meta-prior values throughout the network compared
198 to those used during learning affects imitative behavior. We found that a network that tightly regulates the
199 complexity term by setting smaller values of the meta-prior tends to follow human movement patterns by
200 adapting its internal states. On the other hand, the network that loosely regulates the complexity term by

201 setting larger values of the meta-prior tends to generate more egocentric/self-centered movement patterns
202 with less sensitivity to changes or fluctuations in human movement patterns by adapting its internal state
203 less. The current paper presents a detailed analysis of the underlying mechanisms accounting for these
204 observed phenomena.

205 Below, the *Model* section details the proposed model. It describes an overall system, learning process,
206 derivation of the evidence lower bound of the proposed model, how the trained model was tested in pseudo
207 imitative interaction, and implementation of the model. The *Experiment* section explains the experimental
208 design, procedures of data preparation, and the results of the two experiments. The *Discussion* section
209 summarizes the experimental results and discusses their implications.

2 MODEL

210 2.1 Model overview

211 This subsection describes briefly how multimodal imitative interaction of agents perceiving visuo-
212 proprioceptive sensory inputs can be modeled using concepts of predictive coding and active inference.
213 Among various types of imitation, synchronized imitation is considered in the current study by virtue of
214 its simplicity. In synchronized imitation, the agent is required to imitate target patterns demonstrated by
215 its counterpart by predicting them on the basis of prior learning. Although target patterns to imitate are
216 structurally the same as previously learned patterns, they could involve marginal variations, as in speed,
217 amplitude, and shape. Synchronized imitation can be achieved by means of iterative cycling of sensory
218 input predictions during the demonstration, generation of corresponding movement, and updating the
219 latent state of the network using the resulting sensory prediction error. To generate movement, one step,
220 look-ahead prediction of proprioception is fed into an inverse model (Kawato et al., 1987), which is often
221 implemented by a PID feedback controller in robots. A PID feedback controller computes an optimal
222 motor torque as the motor command to minimize the error between the predicted proprioception (the target
223 joint angles) and the actual proprioception (the actual joint angles). This corresponds to active inference
224 (Friston et al., 2010, 2011), as described previously. The latent state can be updated using a scheme called
225 error regression (Tani and Nolfi, 1999; Ito and Tani, 2004; Hwang et al., 2020; Ahmadi and Tani, 2019), by
226 which sensory perception assumed in a predictive coding framework can be performed.

227 Now we look at how the PV-RNN (Ahmadi and Tani, 2019) can be used to implement the model for
228 multimodal imitative interaction of a robot agent receiving visuo-proprioceptive sensory inputs based on
229 frameworks of predictive coding and active inference. Figure 1 shows the overall system view, consisting
230 of a PV-RNN, a robot, and a human counterpart. The human demonstrates movement patterns to the robot
231 both visually and kinesthetically, guiding the robot's posture via a motion capture suit. Unfortunately, it
232 was infeasible for the proposed system to work stably in real-time because posterior inference using an
233 error-regression scheme, detailed in section 2.4, requires intensive computation. Hence, in the current
234 study, we simulated the imitative interaction between a human and a robot shown in Figure 1 as a pseudo
235 imitative interaction in which pre-recorded body movements sampled from a human serve as the robot
236 counterpart using the setting shown in Figure 1 (C).

237 PV-RNN is considered a generative model, formulated in a continuous spatio-temporal domain, employing
238 a variational Bayes framework, as described previously. It infers the posterior at each time step using
239 variational inference, in which the reconstruction error is minimized with regularization of the KL
240 divergence between the inferred posterior and the conditional prior. This is implemented by means
241 of a so-called *error-regression scheme*, detailed in section 2.4.

242 A PV-RNN inherits the concept of a Multiple Timescale Recurrent Neural Network (MTRNN)(Yamashita
243 and Tani, 2008), which is characterized by its architecture because it allocates different timescale dynamics
244 to different layers. Higher layers are endowed with slower timescale dynamics and lower layers with
245 faster dynamics, as inspired by recent cognitive neuroscience evidence (Newell et al., 2001; Huys et al.,
246 2004; Smith et al., 2006; Kording et al., 2007). Introduction of multiple timescale dynamics can enhance
247 abstraction and generalization in learning by extracting action-primitive hierarchies or chunking structures
248 from observed multimodal sensory inputs (Yamashita and Tani, 2008; Choi and Tani, 2018; Hwang et al.,
249 2020).

250 These characteristics of variational Bayes frameworks and MTRNN enable PV-RNN to utilize
251 hierarchically organized probabilistic representation, i.e., while the network extracts a hierarchical structure
252 from an observation, it also assigns a different degree of uncertainty within the hierarchy. For example,
253 given a task in which the network is required to predict a sequence of body movements comprised of a
254 small number of primitive patterns, the network can be certain about details of the primitive patterns, but
255 less certain about the sequence of the primitives. In such a case, the lower level of the network responsible
256 for prediction of details of each primitive movement shows small uncertainty, while the higher level in
257 charge of prediction of the order of those primitive patterns shows high uncertainty.

258 Sensory modules for proprioception and vision were modeled with multi-layered PV-RNNs and modules
259 were connected via an associative module, also based on a PV-RNN. Figure 1 (A) depicts a schematic of
260 the proposed model and how it is trained. The associative module generates the prior, conditioned by the
261 latent state at the previous time step in this module. The prior is then fed to both the proprioception and
262 vision modules along the top-down pathway. Each sensory module also generates a prior at each time step
263 conditioned by the previous latent state of the module, computed using top-down information provided by
264 the associative module, by which predictions of sensory inputs, proprioception and vision, are generated in
265 the subsequent time step. Note that the vision module predicts a low-dimensional vector, which is then
266 fed to a CNN-type decoder (LeCun et al., 1989, 1998) to generate actual pixel visual images, in order to
267 reduce computational costs.

268 A dataset of visuo-proprioceptive patterns demonstrated by human participants is used to train the model.
269 To generate these data, a human wearing a motion-capture suit demonstrates body movements while
270 simultaneously recording a video. The motion capture suit maps the human's body configuration into the
271 humanoid robot's joint angle values. These synchronized joint-angle trajectories and videos serve as the
272 target of the model. The whole network is optimized simultaneously so as to maximize the evidence lower
273 bound of the model via BPTT. The design of body movement patterns used in this study is detailed in
274 section 3.2.

275 Figure 1 (B) describes how the trained model performs imitative interactions. An imitative interaction
276 involves a cycle of predictions with conditional prior and posterior inferences. At each time step, the
277 network predicts proprioception p_t and a low-dimensional latent representation of vision l_t with the prior
278 conditioned by the latent variable in each module at the previous time step. The proprioceptive prediction
279 p_t is supplied to the controller, followed by computation of motor commands m_t to achieve the expected
280 joint positions and generation of the movement. Then, a new visual image and proprioception are acquired.
281 The raw pixel image is fed to a CNN-type encoder that has been separately trained to obtain the target for
282 the low-dimensional latent representation \bar{l}_t . Resultant prediction errors e_t^l and e_t^p are computed in vision
283 and proprioception, respectively, which are then used to infer the posterior in each PV-RNN layer with
284 regulation of the KL divergence between the inferred posterior and the conditional prior so that the lower
285 bound is maximized by BPTT. This optimization process to infer the posterior is iterated a fixed number

286 of times at each sensory-motor sampling time step, and the optimized posterior is used to make the best
 287 prediction with the conditional prior to the succeeding time step.

288 Figure 1 (C) denotes how the robot's network model senses movement patterns demonstrated by the
 289 human counterpart.

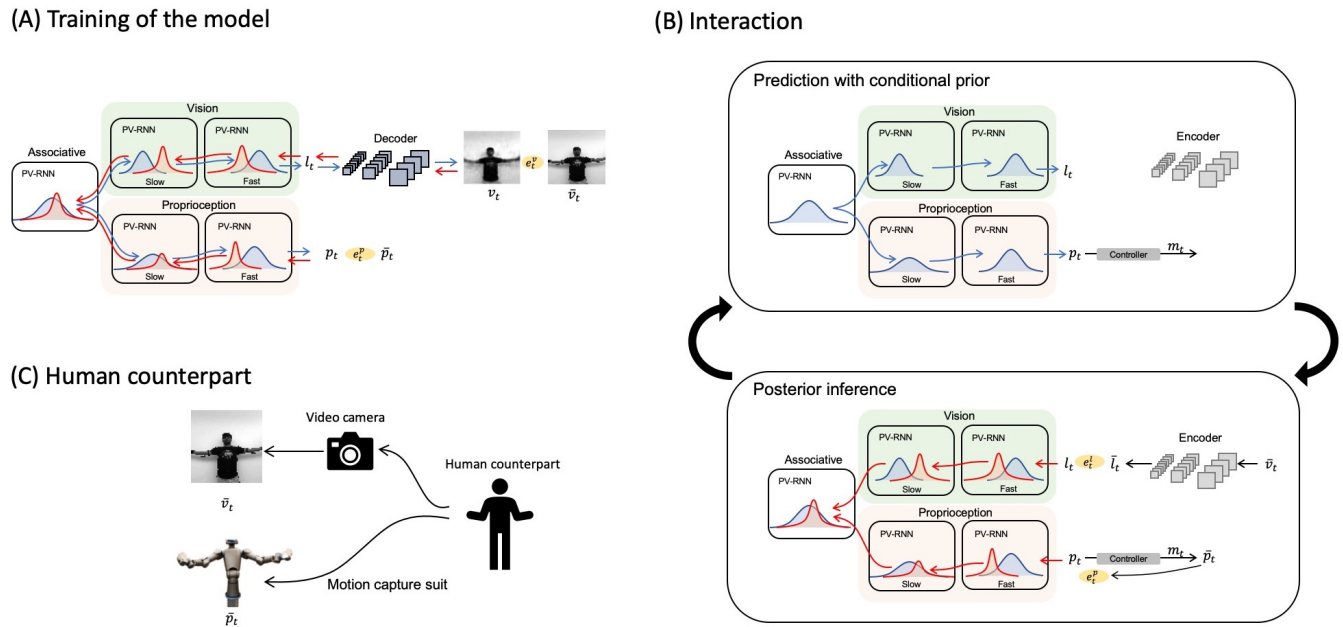


Figure 1: Overall schematic of the proposed model. Blue and red bell curves represent prior and posterior distributions, respectively. Blue and red arrows illustrate information flows of the prediction with conditional prior and posterior inferences, respectively. (A) The training scheme of the proposed network model. (B) The cycle of prediction with conditional prior and posterior inferences during an interaction with a human. (C) A diagram of providing the configuration of a human counterpart to the network.

290 2.2 Derivation of evidence lower bound

291 PV-RNN is a generative, inference model based on the graphical representation shown in Figure 2 (this
 292 figure will be explained in detail in section 2.4). It is comprised of deterministic variables \mathbf{d} , i.e., assumed
 293 to follow Dirac delta distributions, and stochastic variables \mathbf{z} . The model includes a prior and infers the
 294 corresponding posterior by variational inference. We modified the original PV-RNN at four points with
 295 respect to dependencies of variables. First, in our model, there are no connections between the output of
 296 the network \mathbf{x} and \mathbf{z} , which exist in the original PV-RNN. This is for simplification of the model, and it
 297 was confirmed that removing these connections did not hinder network performance. Second, the current
 298 network does not have connections from the lower layer to the higher layer, which the original network
 299 does have. This modification is intended to separate more clearly the information flow between top-down
 300 generative prediction and bottom-up error propagation. Third, diagonal connections from the higher layer
 301 during the previous time step to the lower layer during the succeeding time step are changed to vertical
 302 connections during the same time step. Last, the prior distribution of \mathbf{z}_t at time step 1 has been changed. In
 303 the original study, the distribution is simply mapped from \mathbf{d}_0 . In the current study, however, it is assumed
 304 that $p(\mathbf{z}_1)$ follows a unit Gaussian distribution to control the initial sensitivity of the model. Following
 305 derivation of the evidence lower bound in Ahmadi and Tani (2019) and considering the introduction of
 306 the unit Gaussian at time step 1, the evidence lower bound of the proposed visuo-proprioceptive model is

307 derived as

$$\begin{aligned}
 \ln(\mathbf{p}_{1:T}, \mathbf{v}_{1:T} | \mathbf{d}_0^*) &\geq \sum_{t=1}^T \left\{ \mathbb{E}_{q^a, q^p} [\ln P(\mathbf{p}_t | \mathbf{d}_t^{p,1})] + \mathbb{E}_{q^a, q^v} [\ln P(\mathbf{v}_t | \mathbf{d}_t^{v,1})] \right\} \\
 &\quad - \sum_{l \in A} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p, e_{1:T}^v) \| p(\mathbf{z}^u)] - \sum_{l \in P} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p) \| p(\mathbf{z}^u)] \\
 &\quad - \sum_{l \in V} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^v) \| p(\mathbf{z}^u)] + \sum_{t=2}^T \left\{ - \sum_{l \in A} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] \right. \\
 &\quad \left. - \sum_{l \in P} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] - \sum_{l \in V} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] \right\}
 \end{aligned} \tag{3}$$

308

309 where A, P, and V represent the associative module, the proprioception module, and the vision module,
 310 respectively, and l indicates the index of a layer in each module. $\mathbf{p}_{1:T}$ and $\mathbf{v}_{1:T}$ are time series proprioceptive
 311 and visual patterns. \mathbf{d}_0^* represents \mathbf{d} in all layers at time step 0. \mathbb{E}_{q^a, q^p} denotes the expectation over all
 312 distributions of \mathbf{z} in the associative module and the proprioception module, and \mathbb{E}_{q^a, q^v} denotes expectation
 313 over all distributions of \mathbf{z} in the associative module and the vision module. $\mathbf{d}_t^{p,1}$ is the deterministic variable
 314 in the lowest layer of the proprioception module at time step t , and $\mathbf{d}_t^{v,1}$ is that in the lowest layer of the
 315 vision module. \mathbf{z}_t^l is the stochastic variable at time step t in the l th layer in each module. $e_{t:T}^p$ and $e_{t:T}^v$
 316 are the prediction errors between the predicted patterns and the target patterns at time step from t to T in
 317 proprioception and vision, respectively. $p(\mathbf{z}^u)$ indicates the unit Gaussian distribution serving as the prior
 318 at time step 1. By introducing the meta-prior, which weights the KL divergence between the approximate
 319 posterior and the prior in a layer-specific manner, the evidence lower bound of the model is defined as

$$\begin{aligned}
 \mathcal{L}_w &:= \sum_{t=1}^T \left\{ \underbrace{\mathbb{E}_{q^a, q^p} [\ln P(\mathbf{p}_t | \mathbf{d}_t^{p,1})]}_{\text{Accuracy in proprioception}} + \underbrace{\mathbb{E}_{q^a, q^v} [\ln P(\mathbf{v}_t | \mathbf{d}_t^{v,1})]}_{\text{Accuracy in vision}} \right\} \\
 &\quad - \sum_{l \in A} w_1^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p, e_{1:T}^v) \| p(\mathbf{z}^u)]}_{\text{Complexity in associative module}} - \sum_{l \in P} w_1^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p) \| p(\mathbf{z}^u)]}_{\text{Complexity in proprioception module}} \\
 &\quad - \sum_{l \in V} w_1^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^v) \| p(\mathbf{z}^u)]}_{\text{Complexity in vision module}} + \sum_{t=2}^T \left\{ - \sum_{l \in A} w^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)]}_{\text{Complexity in associative module}} \right. \\
 &\quad \left. - \sum_{l \in P} w^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)]}_{\text{Complexity in proprioception module}} - \sum_{l \in V} w^l \underbrace{D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)]}_{\text{Complexity in vision module}} \right\}
 \end{aligned} \tag{4}$$

320

321 where w_1^l indicates the meta-prior in the l th layer at $t = 1$ in the associative module, the proprioception
 322 module, and the vision module, respectively. w^l represents the meta-priors in the l th layer after $t = 2$ in
 323 each module. Parameters of the model are optimized by maximizing the lower bound, which corresponds
 324 to minimizing the free energy.

325 2.3 Learning process

326 It is noted that unlike other models employing online learning methods (Boucenna et al., 2014, 2016),
 327 our model is trained offline with pre-recorded dataset. The entire network model is trained by maximizing
 328 the evidence lower bound. Thus, given the time step length T of proprioceptive patterns $\mathbf{p}_{1:T}$ and visual

329 patterns $v_{1:T}$, the cost function to be minimized is defined as

$$\begin{aligned}
\text{cost} := & \sum_{t=1}^T \left\{ \frac{1}{2R^p} \|\mathbf{p}_t - \bar{\mathbf{p}}_t\|^2 + \frac{1}{2R^v} \|\mathbf{v}_t - \bar{\mathbf{v}}_t\|^2 \right. \\
& + \sum_{l \in A} \frac{w_1^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p, e_{1:T}^v) \| p(\mathbf{z}^u)] + \sum_{l \in P} \frac{w_1^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^p) \| p(\mathbf{z}^u)] \\
& + \sum_{l \in V} \frac{w_1^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_1^l | \mathbf{d}_0^l, e_{1:T}^v) \| p(\mathbf{z}^u)] + \sum_{t=2}^T \left\{ \sum_{l \in A} \frac{w^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] \right. \\
& \left. + \sum_{l \in P} \frac{w^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^p) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] + \sum_{l \in V} \frac{w^l}{R^l} D_{\text{KL}}[q(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l, e_{t:T}^v) \| p(\mathbf{z}_t^l | \mathbf{d}_{t-1}^l)] \right\}
\end{aligned} \tag{5}$$

331 where A, P, and V represent the associative module, the proprioception module, and the vision module. R^p
332 and R^v are the dimensions of proprioceptive patterns and visual patterns to normalize prediction errors,
333 and R^l is the dimension of the distributions of \mathbf{z} to normalize the KL divergence. Each output in the vision
334 and proprioceptive modules is represented by a multivariate Gaussian distribution with an estimation of
335 the mean for each dimension and covariant matrix as the identity matrix, for simplicity. This leads to
336 minimization of the mean squared error, which is an estimator of the log-likelihood in the accuracy term
337 when maximizing the lower bound.

338 Since the prior and posterior distributions are assumed to follow a multivariate Gaussian distribution with
339 a diagonal covariant matrix, the KL divergence in the cost function is analytically computed. Given two n
340 dimensional multivariate Gaussian distributions $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^p, \boldsymbol{\sigma}^p)$ and $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^q, \boldsymbol{\sigma}^q)$ where
341 $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)^T$ and $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_n)^T$,

$$D_{\text{KL}}[q(\mathbf{z}) \| p(\mathbf{z})] = \sum_{i=1}^n \left\{ \ln \left(\frac{\sigma_i^p}{\sigma_i^q} \right) + \frac{(\mu_i^p - \mu_i^q)^2 + (\sigma_i^q)^2}{2(\sigma_i^p)^2} - \frac{1}{2} \right\} \tag{6}$$

343 The parameters of the model, including an adaptive variable \mathbf{a} introduced in the following section, are
344 optimized using BPTT. To perform error-regression explained in section 2.4, an encoder was also trained
345 separately.

346 2.4 Error-regression with shifting window

347 Ahmadi and Tani (2019) proposed a scheme, the error-regression (ER) with shifting window to test the
348 trained model in a way that is consistent with concepts of predictive coding and active inference. In this
349 scheme, the trained network attempts to predict sensory inputs in the next time step while inferring the
350 posterior in the immediate past window of a fixed length, using the reconstruction error in the window. The
351 window is referred to as the ER window in the following. It is essential to note that ER for maximizing the
352 evidence lower bound is conducted by iterating two processes of forward computation (Figure 2 (A)) and
353 posterior update (Figure 2 (B)) for specific times at each sensory sampling time step.

354 PV-RNN has unique variables \mathbf{a} that facilitate updating the posterior. \mathbf{a} is time step-specific and has
355 the same dimension as \mathbf{z} in each PV-RNN layer. In other words, when a PV-RNN layer with \mathbf{z} with its
356 dimensionality n tries to infer the posterior for the last T time steps inside the ER window, the PV-RNN
357 layer has $n \times T$ \mathbf{a} valuables and updates them to modify the representation of the posterior. A detailed
358 computation scheme of the posterior using the adaptive variable \mathbf{a} is found in section 2.5. Importantly, in
359 ER, weights and biases of the network are fixed, and only the adaptive variables \mathbf{a} are updated.

360 Let us consider an example of error-regression in which the length of the ER window is two time steps,
 361 and the network has two layers, as shown in Figure 2. Figure 2 (A) illustrates the forward computation
 362 at time step t to infer the posterior. In the forward computation, the network computes the conditional
 363 prior, $p(z_{t-1}|\mathbf{d}_{t-2})$ and $p(z_t|\mathbf{d}_{t-1})$, and the posterior, $q(z_{t-1}|\mathbf{d}_{t-2}, e_{t-1:t})$ and $q(z_t|\mathbf{d}_{t-1}, e_t)$ in each
 364 layer, and generates the prediction with sampling from the posterior distribution inside the window. Then,
 365 the reconstruction error e_{t-1} and e_t , and the KL divergence between respective pairs of the conditional
 366 prior and posterior are computed.

367 Based on the reconstruction error and KL divergence, the inferred posterior is updated to maximize the
 368 evidence lower bound. Figure 2 (B) illustrates how the reconstruction error is back-propagated through
 369 variables and layers to \mathbf{a} , which is responsible for updating the posterior. Using the updated posterior, the
 370 network again performs the forward computation. It should be noted that since the $q(z_{t-1}|\mathbf{d}_{t-2}, e_{t-1:t})$ has
 371 been updated, \mathbf{d}_{t-1} is different from the one before the update; thus, $p(\mathbf{z}|\mathbf{d}_{t-1})$ is also changed through the
 372 posterior update. The reconstruction error and the KL divergence are further computed, and the posterior is
 373 updated. This iterative process of forward computation and posterior update is repeated a fixed number
 374 of times to optimize the approximate posterior for maximizing the evidence lower bound computed with
 375 given meta-prior values.

376 After finishing all iterations, the network generates a new sensory prediction \mathbf{x}_{t+1} with a conditional
 377 prior using the inferred posterior inside the ER window. Then the ER window shifts one time step and the
 378 next target sensation \mathbf{x}_{t+1} is sampled, and the forward computation and the posterior update are reiterated
 379 at time step $t + 1$. In the proposed model, the ER is performed for both visual sensation and proprioception
 380 simultaneously, and this scheme of using ER with a shifting window was used to test an imitative interaction
 381 after training the entire network, as will be described later.

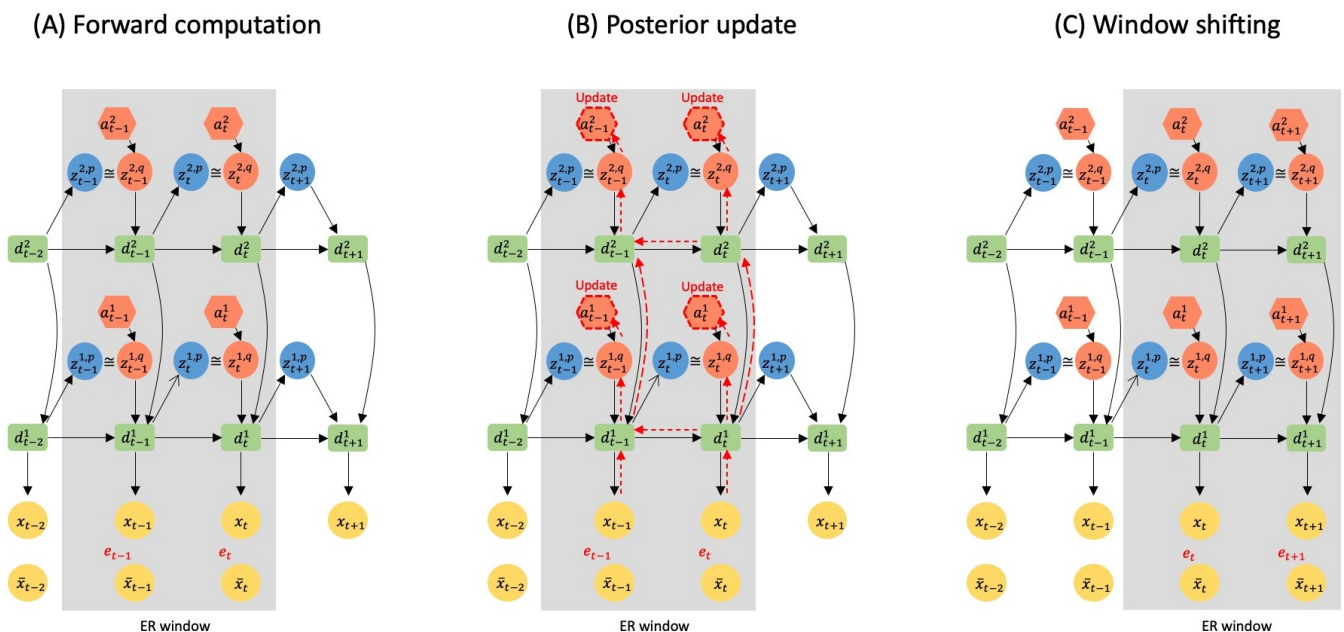


Figure 2: A graphical representation of error-regression with a shifting window. The gray area represents the ER window. Black arrows indicate forward computations. Red arrows indicate how reconstruction errors are propagated to \mathbf{a} inside the ER window by BPTT. (A) illustrates the information flow of forward computation at time step t . (B) shows the update of the posterior inside the ER window at time step t . (C) shows the window shifting to time step $t + 1$.

382 2.5 Model implementation

383 The proposed model for the imitative interaction via visuo-proprioceptive sensation consists of three
 384 modules: an associative module, a proprioception module, and a vision module. This subsection describes
 385 a detailed computation scheme in each module.

386 2.5.1 The associative module

387 The associative module is comprised of a PV-RNN. Since we adopted an MTRNN computation scheme
 388 in PV-RNN, its computations are as follows.

$$389 \mathbf{u}_t^{a,l} = \begin{cases} \mathbf{W}_{dd}^{a,ll} \mathbf{d}_{t-1}^{a,l} + \mathbf{W}_{dz}^{a,ll} \mathbf{z}_t^{a,l} + \mathbf{b}^{a,l} & \text{if top layer} \\ \mathbf{W}_{dd}^{a,ll} \mathbf{d}_{t-1}^{a,l} + \mathbf{W}_{dd}^{a,ll+1} \mathbf{d}_t^{a,l+1} + \mathbf{W}_{dz}^{a,l} \mathbf{z}_t^{a,l} + \mathbf{b}^{a,l} & \text{otherwise} \end{cases} \quad (7)$$

389

$$390 \mathbf{h}_t^{a,l} = \left(1 - \frac{1}{\tau^{a,l}}\right) \mathbf{h}_{t-1}^{a,l} + \frac{1}{\tau^{a,l}} \mathbf{u}_t^{a,l} \quad (9)$$

390

$$391 \mathbf{d}_t^{a,l} = \tanh(\mathbf{h}_t^{a,l}) \quad (10)$$

391

392 where $\mathbf{u}_t^{a,l}$ is the sum of inputs to l th layer of the associative module. $\mathbf{W}_{dd}^{a,ll}$, $\mathbf{W}_{dz}^{a,ll}$, and $\mathbf{W}_{dd}^{a,ll+1}$ are
 393 weight matrices for recurrent connections, the stochastic variable \mathbf{z} , and the input from the higher layer,
 394 respectively. $\mathbf{b}^{a,l}$ is the bias in the l th layer in the associative module, and $\tau^{a,l}$ is the time constant for
 395 MTRNN computation in the l th layer of the associative module. \tanh is the activation function. The
 396 stochastic variable \mathbf{z} is assumed to follow a multivariate Gaussian distribution with a diagonal covariant
 397 matrix, and the deterministic variable \mathbf{d} predicts mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\sigma}$ of the distribution. That is, for
 398 computation of the prior,

$$399 p(\mathbf{z}_t^{p,a,l}) = \begin{cases} \mathcal{N}(\mathbf{z}^u; \mathbf{0}, \mathbf{I}) & \text{if } t=1 \\ p(\mathbf{z}_t^{p,a,l} | \mathbf{d}_{t-1}^{a,l}) = \mathcal{N}(\mathbf{z}_t^{p,a,l}; \boldsymbol{\mu}_t^{p,a,l}, \boldsymbol{\sigma}_t^{p,a,l}) & \text{otherwise} \end{cases} \quad (11)$$

399

$$400 \boldsymbol{\mu}_t^{p,a,l} = \tanh(\mathbf{W}_{\mu d}^{a,l} \mathbf{d}_{t-1}^{a,l} + \mathbf{b}_{\mu}^{a,l}) \quad (13)$$

400

$$401 \boldsymbol{\sigma}_t^{p,a,l} = \exp(\mathbf{W}_{\sigma d}^{a,l} \mathbf{d}_{t-1}^{a,l} + \mathbf{b}_{\sigma}^{a,l}) \quad (14)$$

401

$$402 \mathbf{z}_t^{p,a,l} = \boldsymbol{\mu}_t^{p,a,l} + \boldsymbol{\sigma}_t^{p,a,l} * \boldsymbol{\epsilon} \quad (15)$$

402

403 where $\boldsymbol{\mu}_t^{p,a,l}$ and $\boldsymbol{\sigma}_t^{p,a,l}$ are the mean and variance for the prior distribution of $\mathbf{z}_t^{p,a,l}$ at time step t in l th
 404 layer in the associative module. $\mathbf{W}_{\mu d}^{a,l}$ and $\mathbf{W}_{\sigma d}^{a,l}$ are the weight matrices for $\mathbf{d}_{t-1}^{a,l}$. $\mathbf{b}_{\mu}^{a,l}$ and $\mathbf{b}_{\sigma}^{a,l}$ are the
 405 biases for each computation. \tanh in computation of mean is used for stability of optimization, and \exp in
 406 $\boldsymbol{\sigma}$ is for variance to be positive. $\boldsymbol{\epsilon}$ is sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. To approximate the posterior, PV-RNN has
 407 adaptive variables \mathbf{a} that are specific to time step and sequence. \mathbf{a} is optimized during learning with the
 408 prediction error via BPTT. By considering \mathbf{a} , the computations of posterior are

$$409 q(\mathbf{z}_t^{q,a,l} | \mathbf{d}_{t-1}^{a,l}, \mathbf{e}_{t:T}^p, \mathbf{e}_{t:T}^v) = \mathcal{N}(\mathbf{z}_t^{q,a,l}; \boldsymbol{\mu}_t^{q,a,l}, \boldsymbol{\sigma}_t^{q,a,l}) \quad (16)$$

409

$$410 \boldsymbol{\mu}_t^{q,a,l} = \tanh(\mathbf{W}_{\mu d}^{a,l} \mathbf{d}_{t-1}^{a,l} + \mathbf{a}_{\mu,t}^{a,l} + \mathbf{b}_{\mu}^{a,l}) \quad (17)$$

410

$$411 \boldsymbol{\sigma}_t^{q,a,l} = \exp(\mathbf{W}_{\sigma d}^{a,l} \mathbf{d}_{t-1}^{a,l} + \mathbf{a}_{\sigma,t}^{a,l} + \mathbf{b}_{\sigma}^{a,l}) \quad (18)$$

411

$$412 \mathbf{z}_t^{q,a,l} = \boldsymbol{\mu}_t^{q,a,l} + \boldsymbol{\sigma}_t^{q,a,l} * \boldsymbol{\epsilon} \quad (19)$$

412

413 where $\mu_t^{q,a,l}$ and $\sigma_t^{q,a,l}$ are mean and variance for the posterior distribution of $z_t^{q,a,l}$ at time step t in l th
 414 layer in the associative module. Note that the weight matrices for d are different from those used to
 415 compute the prior. In addition, unlike the peripheral sensory modules of proprioception and vision, the
 416 associative module does not predict the sensory output directly, but rather predicts the latent representation
 417 of visuo-proprioceptive sequences. Therefore, the weights and biases, as well as the adaptive variable a of
 418 the associative module are not optimized instantly from the reconstruction error of the sensory outcomes,
 419 but from the error signals mediated through each sensory module.

420 2.5.2 The proprioception module

421 Proprioceptive patterns are directly generated from the PV-RNN. The highest layer in the proprioception
 422 module receives the input from the lowest layer in the associative layer, and its computations are

$$u_t^{p,l} = \begin{cases} W_{dd}^{pa} d_t^{a,1} + W_{dd}^{p,ll} d_{t-1}^{p,l} + W_{dz}^{p,l} z_t^{p,l} + b^{p,l} & \text{if top layer} \\ W_{dd}^{p,ll} d_{t-1}^{p,l} + W_{dz}^{p,l} z_t^{p,l} + b^{p,l} & \text{otherwise} \end{cases} \quad (20)$$

423

$$424 \quad h_t^{p,l} = \left(1 - \frac{1}{\tau^{p,l}}\right) h_{t-1}^{p,l} + \frac{1}{\tau^{p,l}} u_t^{p,l} \quad (22)$$

$$425 \quad d_t^{p,l} = \tanh\left(h_t^{p,l}\right) \quad (23)$$

426 A proprioceptive pattern at time step t , p_t , is generated from the lowest layer of the proprioception module.

$$427 \quad p_t = \tanh\left(W^p d_t^{p,1} + b^p\right) \quad (24)$$

428 2.5.3 The vision module

429 For the vision module, a scheme to reduce the computation time is introduced. As described in section 2.4
 430 above, in the proposed imitative interaction scheme, the network is required to infer the posterior for the
 431 immediate past at every sensory sampling time step by repeating forward computation and BPTT, which
 432 demands intensive computation. Nevertheless, our model is expected to work in actual robots in real-time in
 433 the future, which necessitates reducing the model's computational complexity. To reduce the computational
 434 demand in the posterior inference in visual perception, we consider a composite network combining a
 435 dynamic PV-RNN and static CNNs for decoding and encoding pixel patterns, instead of introducing full
 436 recurrent connections in this module. In this composite network, when generating predictive output for
 437 the visual input, the PV-RNN part predicts the latent state representation with a relatively low dimension,
 438 which is fed to a CNN decoder to generate the corresponding visual pixel image. On the other hand, when
 439 receiving the visual input, it is transformed to the latent state representation by a CNN encoder. Then, the
 440 prediction error can be computed as the discrepancy in the latent state with a low dimension rather than at
 441 the pixel level with high dimension. This reduces the computational burden significantly for conducting the
 442 BPTT to infer the posterior during imitative interaction. As in the proprioception module, the highest layer
 443 of the vision module receives input from the lowest layer of the associative layer, and its computations are

$$u_t^{v,l} = \begin{cases} W_{dd}^{va} d_t^{a,1} + W_{dd}^{v,ll} d_{t-1}^{v,l} + W_{dz}^{v,l} z_t^{v,l} + b^{v,l} & \text{if top layer} \\ W_{dd}^{v,ll} d_{t-1}^{v,l} + W_{zd}^{v,l} z_t^{v,l} + b^{v,l} & \text{otherwise} \end{cases} \quad (25)$$

444

445

$$\mathbf{h}_t^{v,l} = \left(1 - \frac{1}{\tau^{v,l}}\right) \mathbf{h}_{t-1}^{v,l} + \frac{1}{\tau^{v,l}} \mathbf{u}_t^{v,l} \quad (27)$$

446

$$\mathbf{d}_t^{v,l} = \tanh\left(\mathbf{h}_t^{v,l}\right) \quad (28)$$

447 Then the lowest layer of the PV-RNN predicts the latent state \mathbf{l}_t at time step t , and the visual pattern \mathbf{v}_t is
448 generated by the decoder.

449

$$\mathbf{l}_t = \tanh\left(\mathbf{W}^l \mathbf{d}_t^{v,1} + \mathbf{b}^l\right) \quad (29)$$

450

$$\mathbf{v}_t = \text{decoder}(\mathbf{l}_t) \quad (30)$$

451 In the imitative interaction, the target of latent dynamics $\bar{\mathbf{l}}_t$ of visual patterns $\bar{\mathbf{v}}_t$ at time step t is computed
452 by the encoder.

453

$$\bar{\mathbf{l}}_t = \text{encoder}(\bar{\mathbf{v}}_t) \quad (31)$$

454 To improve the generalization capability of the encoder and decoder, *CoordConv* architecture(Liu et al.,
455 2018) was introduced.

3 EXPERIMENTS

456

3.1 Experimental design

457

458 Using the proposed model, imitative interaction experiments considering human-robot interactions were
459 conducted. Although human-robot interactions ought to be studied in an online fashion to reflect human
460 behavior in response to robot actions, because of the intensive computation required in the error regression
461 scheme, we could not conduct such experiments online. Therefore, the current study examined only
462 the dynamic response of the model network using recorded sequences of visuo-proprioceptive patterns.
463 Therefore, data containing human-demonstrated movement patterns in terms of visuo-proprioceptive
464 sequences were collected both for training the network and for later testing of pseudo-synchronized
465 imitative interaction. After training, the model was tested for pseudo-imitative interaction using novel
466 visuo-proprioceptive patterns with two different scenarios (Experiment 1 and Experiment 2).

466

467 Experiment 1 investigated the issue of coordination and integration of different modalities of sensation
468 by changing the tightness used to regulate the complexity term for each sensory module. For this purpose,
469 the network was trained by assigning different values of the meta-prior to the proprioception and vision
470 modules. We examined the different effects of regulating complexity in the two modules on coordination of
471 different modalities by analyzing them in both the learning process and in the pseudo-imitative interaction
472 tested after learning.

472

473 Experiment 2 investigated the issue of strength of agency as the main motivation of the current study by
474 changing the tightness used to regulate the complexity term for the entire network from that introduced in
475 the training phase. Accordingly, we selected a network trained and evaluated as successful in Experiment 1
476 and then the characteristics of the pseudo-imitative interaction were examined by equally adjusting the
477 meta-priors of each module of this trained network to larger or smaller values.

477

478 In subsequent experiments, some parameters that determine network structure were set as follows. The
479 associative module consisted of a one-layer PV-RNN, and the proprioception module and the vision module
480 consisted of two layers. These PV-RNN layers were characterized by a time scale imposed on MTRNN
481 computation. That is, the higher layer had a larger time constant, producing slow time-scale dynamics,

481 and the lower layer had smaller time constants, generating fast time-scale dynamics. Therefore, in this
 482 study, the higher layer of the proprioception module and the vision module, which receive input from the
 483 associative module, are referred to as the *slow layer*, and the lower layer is referred to as the *fast layer*. As
 484 described in section 2.5.3, the visual perception of the model involves an encoder and a decoder. Their
 485 architectures are summarized in Table 1.

Layer	Kernel size	Stride	Filter	Activation
Encoder				
Conv	33×33	1	5	ReLu
Conv	17×17	1	15	ReLu
Conv	16×16	1	30	tanh
Decoder				
Conv transpose	16×16	1	15	ReLu
Conv transpose	17×17	1	5	ReLu
Conv transpose	33×33	1	1	tanh

Table 1: The architecture of the encoder and the decoder.

486 3.2 Data preparation

487 To obtain a dataset of synchronized visuo-proprioceptive sequences, we used a humanoid robot, Torobo
 488 (Tokyo Robotics Inc.) and a motion capture suit (Perception Neuron, Noitom Ltd.). Torobo is a human-
 489 sized, torso-type humanoid robot with 16 joint-angles, of which 6 are for each arm and 4 are for the torso
 490 and head positions. Human body movements can be mapped to joint-angle trajectories of the robot using
 491 the motion capture suit. A human experimenter wearing the suit demonstrated a set of body movements,
 492 which were mapped as joint-angle trajectories. This demonstration was also recorded with a camera
 493 to obtain corresponding visual patterns. The target sequential movement pattern to be learned by the
 494 robot was designed by considering a probabilistic finite state machine that can generate probabilistic
 495 sequences of three different primitive movement patterns. Those were (A) waving with both arms three
 496 times, (B) rotating the torso to the left with the arms three times, and (C) rotating the torso to the right
 497 with the arms three times. Primitive pattern A is followed either by primitive pattern B or primitive
 498 pattern C with a 50% chance, and primitive patterns B and C are followed by pattern A with a 100%
 499 chance (Figure 3 (A)). One sequence consists of 8 probabilistic transitions of primitive movements. Three
 500 human participants demonstrated and recorded 10 movement sequences each. In other words, the dataset
 501 comprised 30 sequences of visuo-proprioceptive temporal patterns. Recorded visuo-proprioceptive patterns
 502 were down-sampled to 3.75 Hz so that one sequence became 400 time steps. Joint-angle trajectories were
 503 normalized to a range between -1 and 1 . Vision patterns were further converted into gray scale images and
 504 down-sized to 64×64 pixels (Figure 3 (B)). A summary of the training data is shown in Table 2. Visual
 505 trajectories fluctuated far more than proprioceptive trajectories due to various optical conditions, such as
 506 illumination and surface reflectiveness.

	Dimension	time step	Participants	Total sequences
Proprioception	16	400	3	30
Vision	64×64			

Table 2: A summary of the training data.

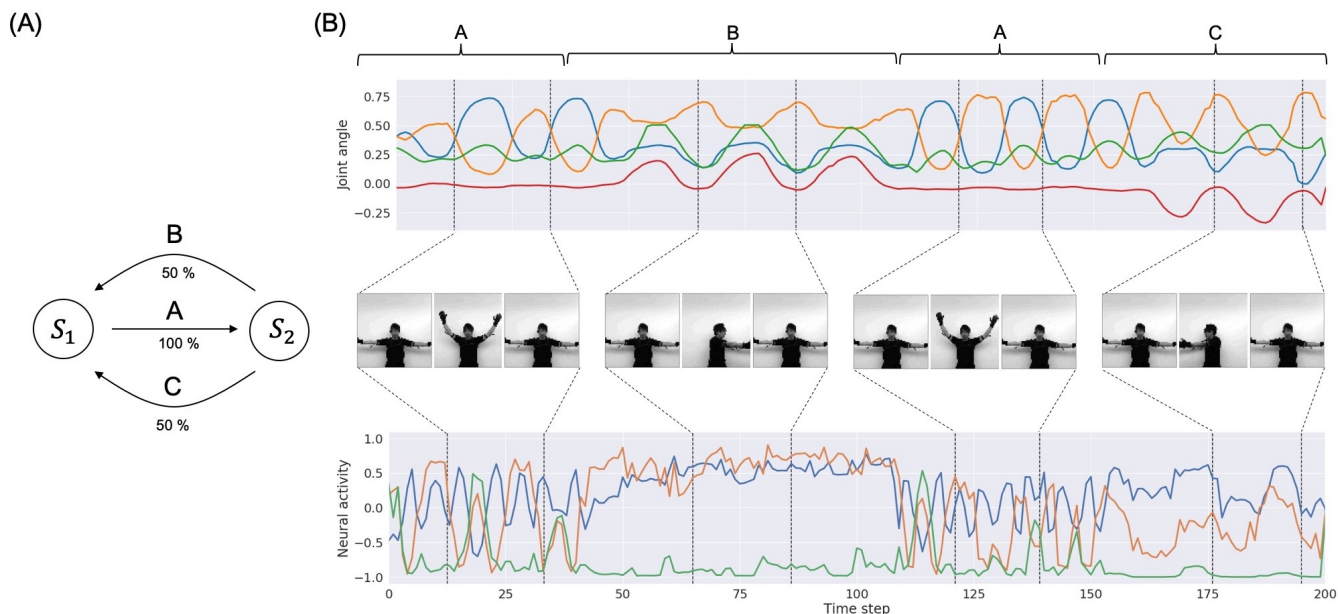


Figure 3: Training data. (A) A diagram of the probabilistic finite state machine. (B) An example of the training dataset. The top row is part of a joint-angle trajectory. The corresponding labels of primitive patterns (A, B, and C) are indicated above the plots. For simplicity, only 4 joint-angles out of 16 representing the movements are shown. The middle row shows corresponding visual pixel images in each period. The bottom row shows visual trajectories in the latent space embedded by the encoder. For simplicity, only three variables out of 20 are shown.

507 Using the training example, the model is required to extract a probabilistic structure such that the primitive
 508 pattern of B or C appears with only a 50% chance after every appearance of the primitive A, by estimating
 509 precision in transitions of primitives with learning. Such learning should be achieved without providing
 510 explicit labels for those primitives, by extracting the underlying chunking and segmentation structure from
 511 continuous sensory signals prepared in the dataset. The PV-RNN can achieve such tasks using a multiple
 512 timescale RNN scheme combined with a Bayesian inference approach (Ahmadi and Tani, 2019).

513 3.3 Experiment 1: Training with different meta-priors in different modalities

514 This experiment investigates effects of changing the tightness used to regulate the complexity term for
 515 each sensory module with regard to coordination and integration of different modalities of sensation. In
 516 addition, this experiment provides successfully trained networks with well-balanced complexity between
 517 the vision and proprioceptive modules for possible use in Experiment 2. To accomplish this, we examined
 518 how assigning different values of the meta-prior to the proprioception and vision modules affects the
 519 learning process and performance in the pseudo-imitative interaction. Two sets of meta-priors w_1 and
 520 w_2 were assigned to the model (Table 3). w_1 has larger values of the meta-prior in the proprioception
 521 module than in the vision module, and they were exchanged in w_2 . Both w_1 and w_2 have the same value
 522 for the meta-prior in the associative module. First, the model was trained with the w_1 and w_2 settings,
 523 and the learning process was examined, with special attention to each component of the lower bound. To

	\mathbb{R}^d	\mathbb{R}^z	τ	w_1 setting		w_2 setting	
				w^l	w_1^l	w^l	w_1^l
Assoc. module	10	1	15	0.0025	0.01	0.0025	0.01
Prop. slow layer	20	2	8	0.005	0.01	0.0025	0.05
Prop. fast layer	30	3	2	0.01	0.01	0.005	0.05
Vision slow layer	20	2	8	0.0025	0.05	0.005	0.01
Vision fast layer	30	3	2	0.005	0.05	0.01	0.01

Table 3: The model configuration in Experiment 1. \mathbb{R}^d and \mathbb{R}^z are the dimensions of d and z , respectively. τ is the time constant of the MTRNN computation in each layer.

524 facilitate training, the Adam optimizer (Kingma and Ba, 2014) was utilized with the parameter settings
 525 $\alpha = 0.001$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. The model was trained 10 times with different random initializations
 526 of model parameters for 10,000 epochs, and the mean and standard deviation of the prediction errors of
 527 proprioception and vision, and the KL divergence of each layer of the model at each epoch were computed.

528 Results are summarized in Figure 4. In comparing w_1 and w_2 conditions, even though the prediction
 529 errors in the proprioception and vision modules showed similar behavior (Figure 4 (A), (B)), the KL
 530 divergence in each module was optimized differently. Despite different values of the meta-prior assigned to
 531 the fast layer of the proprioception module, its KL divergences in w_1 and w_2 conditions were reduced in
 532 exactly the same way (Figure 4 (E)). This is not the case in the fast layer of the vision module (Figure 4 (G)).
 533 The KL divergence in the slow layer of the proprioception module and the slow layer of the vision module
 534 showed different values in w_1 and w_2 settings (Figure 4 (D),(F)). Interestingly, although the associative
 535 module was set to the same value of meta-prior in w_1 and w_2 conditions, the KL divergence in the w_2
 536 setting reached a larger value than in the w_1 setting. This is because the larger value of the meta-prior
 537 assigned to the fast layer of the vision module in the w_2 condition prevented the vision module from
 538 absorbing the fluctuation in observed visual patterns, which resulted in bottom-up fluctuation from the
 539 vision module to the associative module, appearing as a discrepancy between the prior and the posterior in
 540 this module. Because visual sensation contains more inherent randomness than proprioceptive sensation,
 541 as mentioned previously, complexity in this modality should be adequately regulated by setting a smaller
 542 meta-prior value. Otherwise, the discrepancy that appears in the visual module tends to leap to the higher
 543 associative module without being well resolved before.

544 We further tested the trained models in the pseudo-imitative interaction. Training of the models stopped
 545 after 4,000 epochs. Three novel visuo-proprioceptive sequences recorded from three human participants
 546 were prepared for the pseudo-imitative interaction, which also comprised the previous primitive body
 547 movements A, B, and C, the lengths of which were 400 time steps. The length of the ER window was set
 548 to 30 time steps, and the number of optimization iterations for posterior inference by BPTT at each time
 549 step was 30. Namely, at each sensory sampling time step, the network infers the posterior distribution of z
 550 responsible for reconstructing the observation inside the ER window, in which the cycle of the forward
 551 computation and the posterior update described in section 2.4 repeats 30 times. As in learning, Adam was
 552 used to improve optimization with parameter settings $\alpha = 0.2$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. Evaluation of
 553 the error-regression examined how much the reconstruction error in each modality and the KL divergence
 554 at each sub-network in the PV-RNN were minimized. That is, at the point when T' time step window for
 555 the immediate past shifts t time steps, i.e., the current time step is t , the adaptive variable a assigned within
 556 the window is optimized with the iterative process, and at the last iteration, the reconstruction error and the

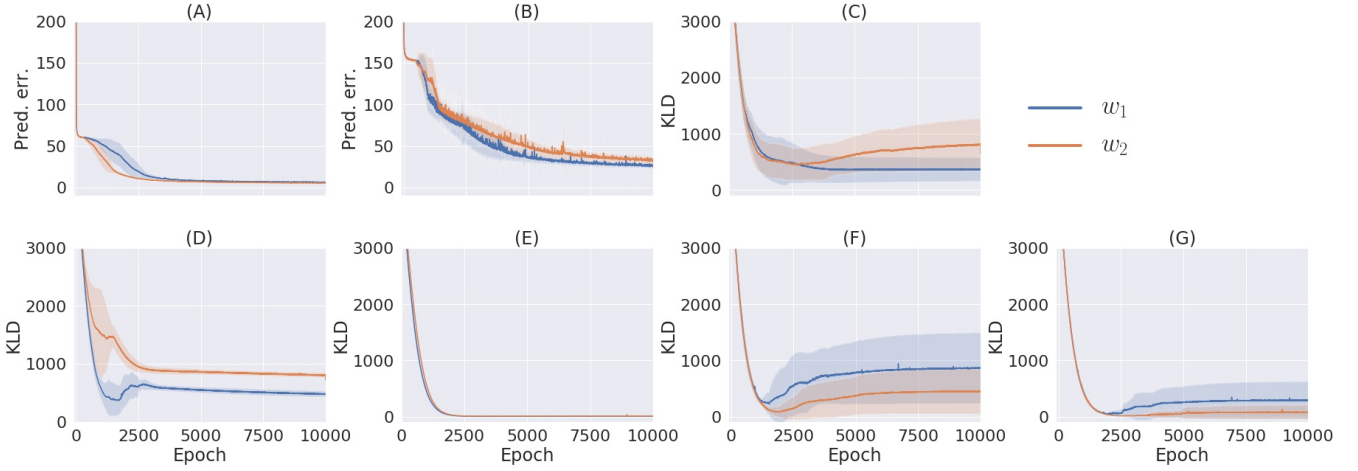


Figure 4: The learning process of the model with two different meta-prior settings. (A) The prediction error in proprioception. (B) The prediction error in vision. (C) The KL divergence in the associative module. (D) The KL divergence in the slow layer of the proprioception module. (E) The KL divergence in the fast layer of the proprioception module. (F) The KL divergence in the slow layer of the vision module. (G) The KL divergence in the fast layer of the vision module. The shadows are the standard deviation of 10 trials with different parameter initializations. Note that values of prediction errors are the sum of the prediction errors at all time steps and sequences normalized by the data dimension.

557 KL divergence are computed inside the window. Therefore, they are defined as

$$558 \quad \text{Proprioception error} := \frac{1}{T} \sum_{t=1}^T \frac{1}{T'} \sum_{t'=1}^{T'} \frac{1}{R^p} \|\mathbf{p}_{t'} - \bar{\mathbf{p}}_{t'}\|^2 \quad (32)$$

$$559 \quad \text{Vision error} := \frac{1}{T} \sum_{t=1}^T \frac{1}{T'} \sum_{t'=1}^{T'} \frac{1}{R^v} \|\mathbf{l}_{t'} - \bar{\mathbf{l}}_{t'}\|^2 \quad (33)$$

$$560 \quad \text{KLD} := \frac{1}{T} \sum_{t=1}^T \frac{1}{T'} \sum_{t'=1}^{T'} \frac{1}{R^z} D_{\text{KL}}[q(\mathbf{z}_{t'} | \mathbf{d}_{t'-1}, e_{t':T}) \| p(\mathbf{z}_{t'} | \mathbf{d}_{t'-1})] \quad (34)$$

561 where t' is the time step inside the window. R^p and R^l are the dimension of proprioception and the latent
 562 space of vision, respectively. R^z is the dimension of \mathbf{z} , and the KL divergence is computed for every
 563 PV-RNN submodule. Models trained in previous experiments were used. The pseudo-imitative interaction
 564 experiment was run 10 times with different random number seeds, and the mean and standard deviation of
 565 each quantity were computed. In addition, one-step, look-ahead prediction error, the discrepancy between
 566 the prediction in the next time step of the current window and the observation, was computed in the vision
 567 module to evaluate prediction accuracy.

568 Figure 5 exemplifies how the pseudo-imitative interaction developed in the $w1$ setting in time-lapse. For
 569 clarity, only parts involving the proprioceptive interaction are shown. Each column shows the representation
 570 of the network when the network finished a posterior inference and made a new prediction at each time
 571 step. The first, second, and third row show representations in the associative module, the slow layer in the
 572 proprioception module and the fast layer in the proprioception module, respectively. Solid lines indicate
 573 the activity of three randomly chosen \mathbf{d} neurons, and dashed lines indicate the KL divergence value at each
 574 time step in each layer. The fourth row shows joint-angle trajectories. Solid lines are predictions generated
 575 by the network, and dashed lines are joint-angle values demonstrated by the human counterpart in the
 576 recorded data. The bottom row shows the reconstruction error, inside the ER window, which was minimized

577 by updating α via BPTT under regularization by the KL divergence between the inferred posterior and the
 578 conditional prior. In section 2.4, describing the error-regression scheme, the network is illustrated in a way
 579 that it only makes the prediction at next time step during the interaction. In this experiment, however, the
 580 network was allowed to generate the prediction not only at next time step, but also at subsequent time steps
 581 with the conditional prior to visualize the network's long-term prediction. This is also the case in Figure 8.

582 At each time step, the network receives a new sensation, computes the reconstruction error and the KL
 583 divergence within the ER window, updates the α such that the lower bound inside the ER window is
 584 maximized, and modifies the prediction after the current time step with the conditional prior. In Figure
 585 5, the network continually modified the future prediction as a result of the posterior inference. Since the
 586 lower bound summed over time steps inside the ER window is maximized, all α s inside the ER window
 587 are updated so that the sum of the reconstruction error and the KL divergence weighted by the meta-prior
 588 inside the ER window are minimized. Therefore, it is often observed that the value of the reconstruction
 589 error or the KL divergence at a certain time step inside the ER window becomes larger at the next time
 590 step, which is considered a transient process in the optimization wherein the past is re-interpreted and
 591 re-represented in coping with a new entering sensation, in terms of post-diction (Shimojo, 2014). In the w_1
 592 setting, larger values of the meta-prior are assigned to lower layers of the network and smaller to higher.
 593 In other words, KL divergences in lower layers are weighted more in the lower bound, and while those
 594 in higher layers are weighted less. Therefore, KL divergences in lower layers were reduced more, and
 595 those in higher layers remained larger after iterative optimization. Owing to MTRNN characteristics of
 596 different time-scales among layers, higher layers showed slower dynamics and lower layers showed faster
 597 dynamics. It is assumed that higher levels predict switching of primitives and lower levels predict sensory
 598 profile changes at each time step. Detailed analysis of this issue was not conducted in the current study
 599 since similar phenomena using MTRNN have been reported frequently (e.g. Yamashita and Tani (2008);
 600 Hwang et al. (2020)).

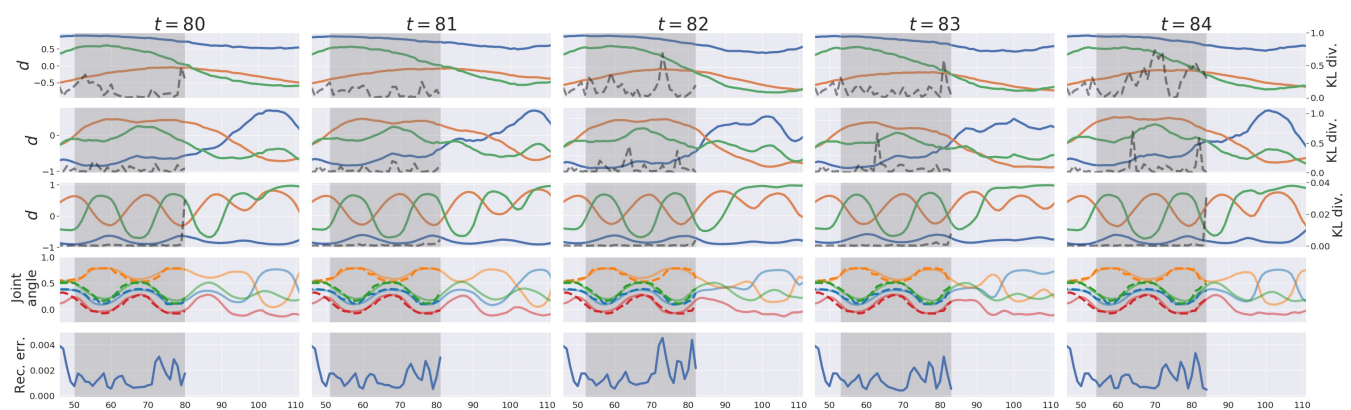


Figure 5: An example of the network representation during testing in w_1 setting. Gray areas indicate the ER window. The first, second, and third rows show representations in the associative module, the slow layer in the proprioception module, and the fast layer in the proprioception module, respectively. Solid lines represent activities of three randomly chosen d neurons, and dashed lines represent the value of the KL divergence at each time step. The 4th row shows predictions (solid lines) and sensations (dashed lines) of joint-angles. For clarity, only four joint-angles of 16 are shown. The bottom row shows the reconstruction error in proprioception.

601 Experimental results are summarized in Table 4. The reconstruction error in proprioception was
 602 remarkably minimized compared to that in vision, in both conditions w_1 and w_2 . This is because vision

603 involves more noise than proprioception. The reconstruction error in vision was smaller for the w_1
 604 condition than the w_2 condition. Furthermore, the KL divergence in the associative module was reduced
 605 more significantly in the w_1 condition than the w_2 condition. This occurred because the vision module
 606 generalized better with noisy visual patterns in the test of pseudo-imitative interaction in the w_1 case than
 607 in the w_2 case by minimizing the complexity term more. Because fluctuation or randomness in visual
 608 sensation was well resolved in the vision module in the w_1 case, the associative module became relatively
 609 free from such fluctuation, as evidenced by the smaller KL divergence observed in the associative module.
 610 As a result, the one-step, look-ahead prediction was also more accurate.

	Proprioception reconstruction error	Vision reconstruction error	Associative KLD	Proprioception slow KLD
w_1	$0.017 \pm 9.5 \times 10^{-4}$	$0.12 \pm 6.9 \times 10^{-3}$	2.0 ± 0.091	1.6 ± 0.12
w_2	$0.011 \pm 2.9 \times 10^{-4}$	$0.19 \pm 1.5 \times 10^{-2}$	3.2 ± 0.15	2.6 ± 0.11

	Proprioception fast KLD	Vision slow KLD	Vision fast KLD	Vision one-step prediction error
w_1	0.57 ± 0.040	8.1 ± 0.17	0.50 ± 0.048	0.20 ± 0.0097
w_2	0.57 ± 0.041	1.9 ± 0.071	0.54 ± 0.048	0.24 ± 0.017

Table 4: The result of the pseudo-imitative interaction experiment. The errors are the standard deviation of 10 different trials with different random number seeds.

611 3.4 Experiment 2: Imitation with stronger and weaker agency

612 This experiment was devised to reveal possible effects of changing the tightness used to regulate
 613 the complexity term for the entire network on the strength of agency exerted in imitative interaction.
 614 Accordingly, we investigated how changes of meta-prior values of the entire network from default values
 615 used in learning affect performance characteristics in the pseudo-imitative interaction. We used a network
 616 that was trained for 4,000 epochs in Experiment 1 with the w_1 setting as the default network. Five meta-prior
 617 settings were prepared for testing of imitative interaction: from smaller values of the meta-prior setting W_1
 618 to the larger setting W_5 with a consistent ratio among all layers of all modules (Table 5). Imitative interaction
 619 with different meta-prior settings was performed with the novel visuo-proprioceptive patterns used in
 620 Experiment 1. Interactions were analyzed in terms of the quantities introduced in previous experiments.
 621 In addition, one-step look-ahead prediction error in proprioception was also measured. Each test with
 622 a different meta-prior setting was repeated with 10 network models trained with different initialization
 623 weights, but with the same parameters for the purpose of examining these quantities statistically.

	W_1	W_2	W_3	W_4	W_5
Associative module	2.5×10^{-5}	2.5×10^{-4}	2.5×10^{-3}	2.5×10^{-2}	2.5×10^{-1}
Proprioception slow layer	5.0×10^{-5}	5.0×10^{-4}	5.0×10^{-3}	5.0×10^{-2}	5.0×10^{-1}
Proprioception fast layer	1.0×10^{-4}	1.0×10^{-3}	1.0×10^{-2}	1.0×10^{-1}	1.0
Vision slow layer	2.5×10^{-5}	2.5×10^{-4}	2.5×10^{-3}	2.5×10^{-2}	2.5×10^{-1}
Vision fast layer	5.0×10^{-5}	5.0×10^{-4}	5.0×10^{-3}	5.0×10^{-2}	5.0×10^{-1}

Table 5: The values of meta-prior in Experiment 2.

624 Results are summarized in Figure 6. As a whole, with smaller values of the meta-prior, the reconstruction
 625 error was minimized more (Figure 6(A)), and the KL divergence remained large (Figure 6(D)), whereas
 626 with larger values of the meta-prior, the KL divergence was minimized more (Figure 6(D)), and the
 627 reconstruction error remained large (Figure 6(A)). This tendency can also be seen in the local proprioception
 628 module and vision module, although the reconstruction error in the vision module was not significantly
 629 different. In the proprioception module, as values of the meta-prior increased, the reconstruction error in
 630 proprioception became large (Figure 6(B)), and the KL divergence became small, both in the slow layer
 631 (Figure 6(F)) and in the fast layer (Figure 6(G)). In the vision module, as values of the meta-prior increased,
 632 though the reconstruction error in vision did not increase as significantly (Figure 6(C), the KL divergence
 633 became small in both the slow layer (Figure 6(H)) and the fast layer (Figure 6(I)). The KL divergence
 634 in the associative module also increased as values of the meta-prior increased (Figure 6(E)). In addition,
 635 with smaller values of the meta-prior, the one-step, look-ahead prediction error was minimized in both
 636 proprioception (Figure 6(J)) and in vision (Figure 6(K)).

637 This is because the KL divergence term in the evidence lower bound was weighted more for minimization
 638 than was the reconstruction error term. In this situation, the posterior $q(\mathbf{z}_t|\mathbf{d}_{t-1}, e_{t:T'})$ at each time step in
 639 the ER window approached its prior $p(\mathbf{z}_t|\mathbf{d}_{t-1})$ by modulating the adaptive value \mathbf{a}_t , which is fed into
 640 the computation of the posterior $q(\mathbf{z}_t|\mathbf{d}_{t-1}, e_{t:T'})$, while the prior $p(\mathbf{z}_t|\mathbf{d}_{t-1})$ was less changed. This means
 641 that network dynamics were driven mainly by the prior, and were less affected by sensory inputs. Network
 642 dynamics become more egocentric by following the prior, which was less modified by looser regulation of
 643 the complexity term (i.e., more weighting for the KL divergence term). On the other hand, with tighter
 644 regulation (i.e., less weighting of the KL divergence term), network dynamics became more adaptive
 645 to changes or fluctuations of sensory inputs by freely modulating the posterior in the direction of error
 646 minimization without being much constrained by the prior. In this condition, the prior $p(\mathbf{z}_t|\mathbf{d}_{t-1})$ at each
 647 time step in the window also changes because the posterior $q(\mathbf{z}_{t-1}|\mathbf{d}_{t-2})$ at the previous time step, which
 648 is mapped to $p(\mathbf{z}_t|\mathbf{d}_{t-1})$ through \mathbf{d}_t also changes.

649 In the course of pseudo-imitative interaction, when the network observes a single time step of a new
 650 sensation, it infers sequences of the posterior inside the ER window with the aforementioned iterative
 651 computation of the error regression. Figure 7 displays some examples of the posterior inference during the
 652 process in which tight regulation of the complexity term (W_1 setting) (Figure 7 (A)) and loose regulation of
 653 the complexity term (W_5 setting) (Figure 7 (B)) are compared. For clarity, part of the network responsible
 654 for proprioception is shown. The columns illustrate, given a single time step of sensory observation, how
 655 the network inferred the posterior in terms of parameters of the posterior distribution, mean $\boldsymbol{\mu}$, and variance
 656 $\boldsymbol{\sigma}$ of multivariate Gaussian distributions under the effect of different values of the meta-prior through
 657 iterations. From the left, each column shows network dynamics before the inference, after 5th, 10th, 15th,
 658 20th, 25th, and 30th iteration of the update of adaptive variable \mathbf{a} inside the ER window with BPTT. The
 659 first, third, and fifth rows plot the relationship among the mean of the prior $\boldsymbol{\mu}^p$ (blue lines), the mean of the
 660 inferred posterior $\boldsymbol{\mu}^q$ (red lines), and the KL divergence (dashed black lines) in the associative module, in
 661 the slow and fast layers of the proprioception module, respectively. The second, fourth, and sixth rows plot
 662 the variance of the prior $\boldsymbol{\sigma}^p$ (blue lines), the variance of the inferred posterior $\boldsymbol{\sigma}^q$ (red lines), and the KL
 663 divergence (dashed black lines) in the associative module, in the slow and fast layers of the proprioception
 664 module, respectively. Although dimensions of \mathbf{z} in the fast layer and in the slow layer of the proprioception
 665 module are greater than one, only one dimension is plotted for visibility.

666 In W_1 setting, the network is assigned smaller values of the meta-prior, which means that the complexity
 667 term is tightly regulated. Therefore, during the course of posterior inference, the inferred posterior is

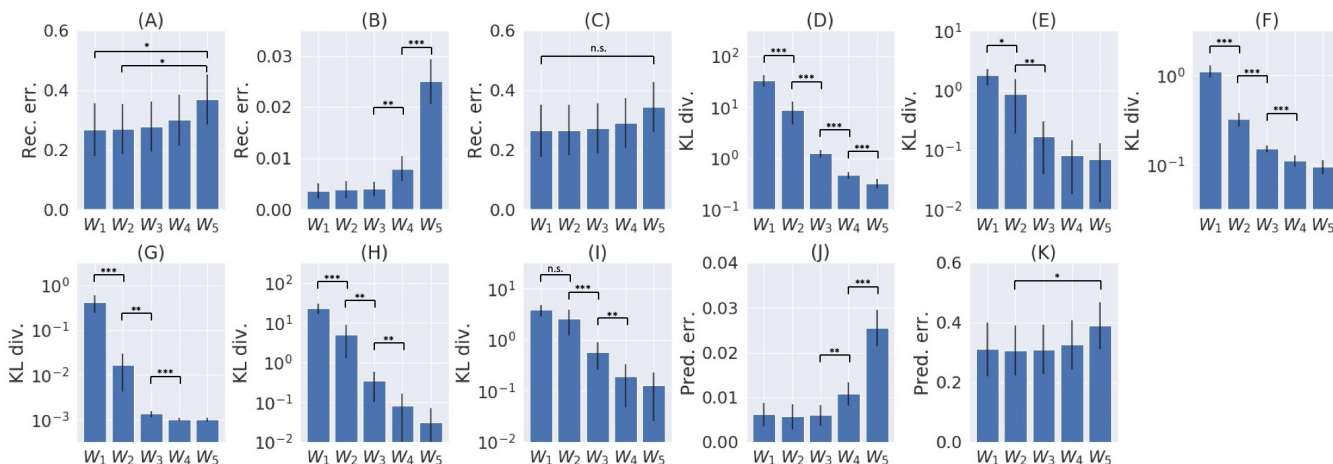


Figure 6: Reconstruction error, KL divergence minimization, and one-step, look-ahead prediction error in error-regression with five meta-prior settings. (A) Sum of reconstruction errors in proprioception and vision. (B) The reconstruction error in proprioception. (C) The reconstruction error in vision. (D) Sum of the KL divergence in all layers. (E) The KL divergence in the associative module. (F) The KL divergence in the slow layer of the proprioception module. (G) The KL divergence in the fast layer of the proprioception module. (H) The KL divergence in the slow layer of the vision module. (I) The KL divergence in the fast layer of the vision module. (J) One-step, look-ahead prediction error in proprioception. (K) One-step, look-ahead prediction error in vision. Error bars represent the standard deviation of 10 models with different weight initialization. Asterisks represent the statistical significance in t-tests: * for $p < 0.05$, ** for $p < 0.01$, and *** for $p < 0.001$. Note that each graph has a different scale.

668 allowed to deviate somewhat from the prior to minimize the reconstruction error compared to the W_5
 669 setting with looser regulation. This can be seen in Figure 7 (A). In the leftmost column, the network
 670 encountered a large reconstruction error in the last time step inside the ER window. This reconstruction
 671 error was eventually resolved while updating the posterior repeatedly as a result of distributing the KL
 672 divergence over the entire network in consideration of values of the meta-prior assigned to each layer.
 673 In the W_1 setting, the associative module had the smallest value of the meta-prior, the slow layer of the
 674 proprioception module had one with a moderate value, and the fast layer of the proprioception module
 675 had the largest value of the meta-prior. Thus, the largest discrepancy between the inferred posterior and the prior
 676 was allowed in the associative module and the smallest discrepancy in the fast layer of the proprioception
 677 module. This can be confirmed by comparing the posterior, the prior, and the value of KL divergence in
 678 each layer in Figure 7 (A).

679 In contrast, in the W_5 setting, the complexity term is loosely regulated with larger values of the meta-prior,
 680 which forces the network to keep the KL divergence small during the posterior inference. This can be
 681 observed in Figure 7 (B). During the iteration, the value of the KL divergence was strongly suppressed, and
 682 as a result, the reconstruction error remained large even after the posterior update. Compared to Figure 7
 683 (A), the posterior was inferred so that it was closer to the prior (red lines are closer to blue lines).

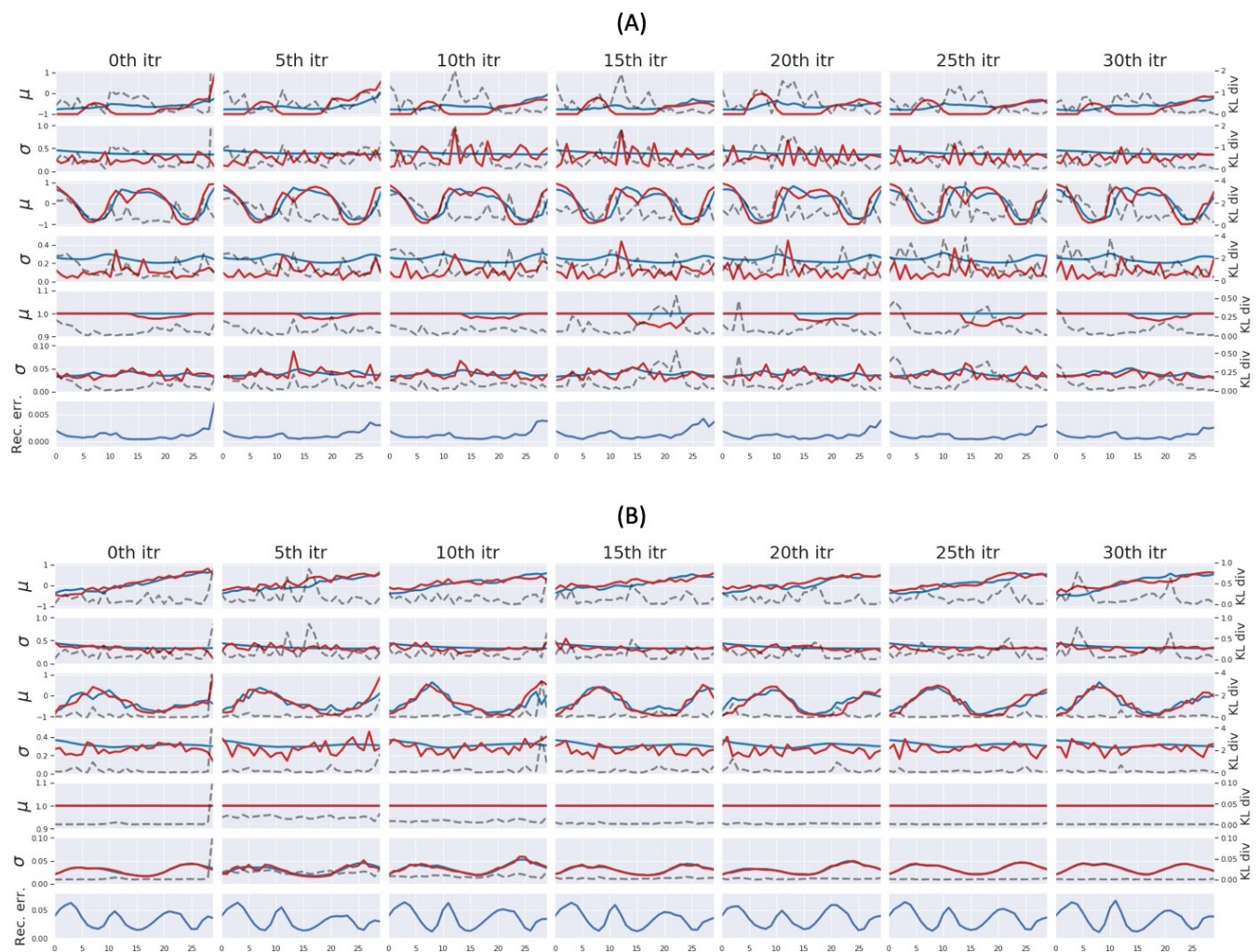


Figure 7: An example of the posterior inference during the pseudo imitative interaction in the W_1 setting (A) and in the W_5 setting (B). For clarity, only those parts involved in the proprioception module are shown. From the left, each column represents the network representation inside the ER window before the inference, after every 5th iteration up to the 30th iteration of the posterior inference. The first, third, and fifth rows show time trajectories of the mean μ of the z in the associative module, the slow layer of the proprioception module, and the fast layer of the proprioception module, respectively. The blue and red lines represent the prior μ^p and the inferred posterior μ^q , respectively. The second, fourth, and sixth rows show the time trajectories of variance σ of z in the associative module, the slow layer of the proprioception module, and the fast layer of the proprioception module, respectively. The blue and red lines indicate the prior σ^p and the inferred posterior σ^q , respectively. Dashed black lines indicate values of the KL divergence in each layer. The bottom row shows the reconstruction error at corresponding time steps.

684 Figure 8 (A) and (B) show examples of time-series plots of related neural activities of the proprioception
 685 module, comparing cases of tight (W_1 setting) and loose (W_5 setting) regulation of the complexity term.
 686 Both cases are computed for a situation observing the same visuo-proprioceptive sequence pattern. With
 687 tight regulation of the complexity term (Figure 8 (A) top), the observation of the primitive A (dashed lines)
 688 was well reconstructed (solid lines) inside the ER window (gray area) from time steps 120 to 150, due to
 689 relatively stronger weighting of the accuracy term compared to the W_5 setting. Plots after time step 150
 690 represent future predictions of the expectation of encountering the primitive B. From time steps 150 to
 691 180 (Figure 8 (A) bottom), the agent observed new sensory information where the primitive C instead of
 692 the predicted primitive of B was encountered. (Remember that there is a 50% chance of encountering the

693 primitive B or the primitive C.) This new observation was reconstructed inside the ER window. Based on
 694 the inferred posterior during this period, the robot updated the future prediction after time step 180 as the
 695 primitive C to be continued. Because of relatively stronger weighting in the accuracy term, the posterior
 696 was inferred to adapt to reality. The prediction was also updated accordingly (Figure 8 (A) bottom).

697 In the case of loose regulation (Figure 8 (B) top), the observation was still well reconstructed inside
 698 the ER window. This is because primitive pattern A always follows either primitive pattern B or C so
 699 that it is easy to predict primitive A. Therefore, the reconstruction error inside the ER window was small
 700 from the beginning. Plots after time step 150 represent future predictions expecting primitive pattern B
 701 to be encountered. After observing new sensory information in which primitive pattern C instead of the
 702 predicted primitive pattern B was encountered between time steps 150 and 180 (Figure 8 (B) bottom);
 703 however, the new observation was not reconstructed well inside the ER window. Due to tight regulation of
 704 the KL divergence term (loose regulation of the complexity term), the posterior was forced closer to the
 705 prior by ignoring the new observation. Consequently, the inferred posterior did not affect the prior as much
 706 as in the W_1 setting, which resulted in generation of consistent predictions for the future. Actually, the
 707 look-ahead prediction made at time step 150, shown in the top row, and the one made at time step 180 in
 708 the bottom row are almost the same. These observations imply that both the prediction of the future and the
 709 reflection of the past become more adaptive to sensory observation in the case of tighter regulation of the
 710 complexity term, whereas they become more persistent regardless of sensory observations in the case of
 711 looser regulation.

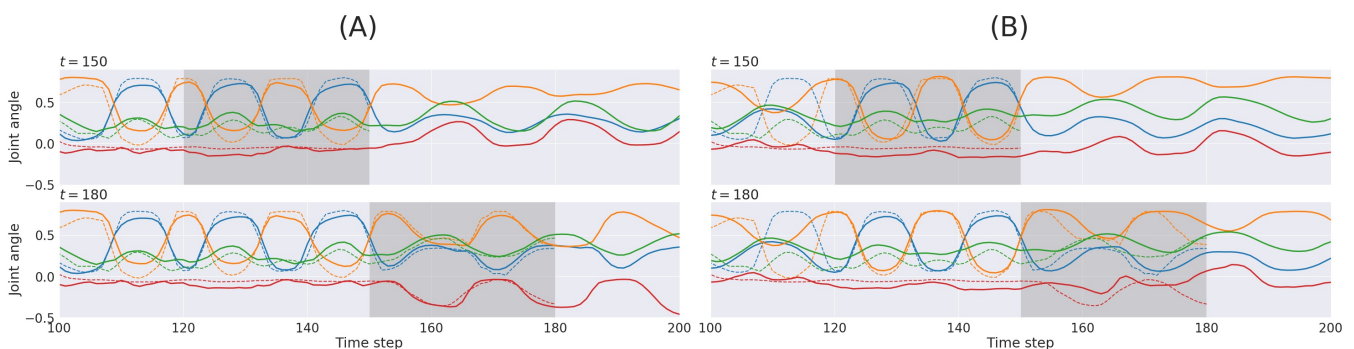


Figure 8: An example of time-series plots of neural activities in the output layer of the proprioception module in the W_1 setting (A) and in the W_5 setting (B). Reconstruction of the past observation and the future prediction at time step 150 (top) and at time step 180 (bottom) are shown. Solid lines represent prediction outputs, and dashed lines represent observations. The shadowed area indicates the error-regression window. For simplicity, only 4 of 16 joint-angles representing movements are shown.

712 Some representative videos related to Experiment 2 can be seen at video link A and at video link B
 713 for the W_1 condition and the W_5 condition, respectively. These videos show how prediction of the future
 714 as well as reflection of the past can be performed for each condition. Also, further temporal details
 715 during the error-regression process can be seen at video link C and at video link D for the W_1 condition
 716 and the W_5 condition, respectively. In these videos, there is some divergence between the prior and the
 717 posterior in terms of mean and variance and they are dynamically changing inside the ER-W in the W_1
 718 condition, whereas these two profiles approximate each other, showing relatively persistent patterns in the
 719 W_5 condition. These observations accord with our analysis, described previously.

4 DISCUSSION

720 The current study investigated underlying mechanism of the strength of agency in social interaction by
721 proposing a model for imitative interaction using multimodal sensation based on the framework of PC and
722 AIF. We proposed a hypothesis that tightness used to regulate the complexity term in the evidence lower
723 bound in the proposed model should contribute to the strength of agency. This hypothesis was evaluated by
724 conducting simulation experiments on a pseudo-human-robot imitative interaction using the model.

725 First, we examined possible effects of changing the tightness used to regulate the complexity term for
726 each sensory module during the learning phase in coordination and integration of different modalities
727 of sensation and those in the test imitation phase. Our results showed that the complexity term in the
728 vision module should be regulated more than that of the proprioception module. This is because vision
729 and proprioception are significantly different with respect to their intrinsic randomness, as visual inputs
730 fluctuate more due to optical conditions, such as illumination and surface reflectiveness. We concluded that
731 the complexity term in the vision module should be regulated much more than that for the proprioception
732 module to achieve better generalization in learning.

733 Next, we investigated the strength of agency as the main focus of the current study by changing the
734 tightness used to regulate the complexity term for the entire network relative to that which was introduced
735 in the training phase. For this purpose, characteristics of pseudo-imitative interaction were examined by
736 scaling the meta-prior of each module equally to larger or smaller values using the network that had been
737 evaluated as successful in the previous experiment. Our results demonstrated that changing the meta-prior
738 this way affects performance characteristics of imitative interaction significantly. With looser regulation of
739 the complexity term, the agent tends to act more egocentrically, without adapting to the other. In contrast,
740 with tighter regulation of the complexity term, the agent tends to follow its human counterpart by adapting
741 its internal state. This result implies that the strength of SoA can be modulated by adjusting the tightness
742 with which the complexity term is regulated after the learning phase.

743 In the current study, we evaluated this hypothesis by considering a task of imitative interaction between a
744 robot and a human counterpart. In such an imitative interaction, there could be two situations: the robot
745 follows the human's movements, or the human follows the robot's movements. In our experimental results,
746 the agent with tight regulation of the complexity term corresponded to the former case, and that with loose
747 regulation to the latter. These findings could provide new insights into computational modeling studies of
748 MNS. Our group's previous studies (Ahmadi and Tani, 2017; Hwang et al., 2020) on modeling MNS using
749 deterministic RNNs that were applied to robot imitation experiments, introduced a scheme similar to the
750 ER scheme described in the current study, in the sense that both reinterpret past observations and update
751 future predictions. In the model, deterministic latent variables at the onset time step of the immediate
752 past window are updated by means of the ER scheme. Since these latent variables are not constrained by
753 any prior probability distribution (unlike the sequential prior scheme), they adapt to sensory sequences
754 encountered for minimizing the error directly wherein the speed of updating is simply determined by the
755 adaptation rate to update the latent variables.

756 On the other hand, in the case of the ER, which uses PV-RNN, the update of stochastic latent variables z
757 at each time step inside the ER window are constrained by the sequential prior represented in terms of a
758 Gaussian probability distribution. If the PV-RNN is developed more toward deterministic dynamics by
759 setting the meta-prior with larger values, the sequential prior for each stochastic latent variable should
760 have a peaky distribution with relatively small variance. In such a case, the approximate posterior cannot
761 adapt to the sensory sequence by using the propagated error signal because the current prior is estimated
762 with a strong belief. In contrast, if the PV-RNN is developed toward a more random process by setting the

763 meta-prior with smaller values, at each time step the prior should exhibit a wide distribution with large
764 variance. Then, the posterior can easily adapt to the sensation using the error signal, because the current
765 prior is estimated with a weak belief. Therefore, the PV-RNN can show both mirror neuron-type adaptive
766 response and egocentric behavior, depending on the setting of the meta-prior in interactions among agents.
767 The deterministic RNN models shown in Ahmadi and Tani (2017); Hwang et al. (2020), however, can only
768 show mirror neuron-like adaptive responses.

769 By following the above discussion, one essential advantage of using variational RNNs, such as PV-RNN,
770 compared with conventional deterministic RNNs, is that they can predict not only future contents, but can
771 also estimate predictability of such predictions or in other words, the credibility of prediction, as discussed
772 in formulation of the *free-energy principle* (Friston, 2005). This sort of cognitive competency of second
773 order prediction by means of representing the belief of prediction, by which the strength of agency can be
774 mechanized, provides modeling of agents, including cognitive robots with more complexity and richness in
775 ways of interacting with other agents, as well as the physical world, as the current study demonstrates, at
776 least partially.

777 In everyday social interactions, humans don't just follow others, nor do they lead them all the time. Rather
778 humans sometimes follow others and sometimes lead them, depending on the moment-by-moment context
779 or social situation. Psychological studies indicate that turn-taking between following and leading can occur
780 quite spontaneously in various social cognitive behaviors, including conversation (Sacks et al., 1978),
781 mother-infant pre-verbal communication (Trevarthen, 1979) and imitation (Nadel, 2002). In considering
782 possible mechanisms underlying turn-taking, some researchers (Ikegami and Iizuka, 2007; Ito and Tani,
783 2004) suggest that turn-taking may develop due to potential instability, such as chaos formed in coupled
784 dynamics between two agents in their modeling studies. We consider meta-level dynamics coupling two
785 agents, whereby the value of the meta-prior to regulate the complexity terms in the two agents counteract
786 one another mutually. This could result in autonomous shifts between the leading mode by increasing the
787 meta-prior and the following mode by reducing it.

788 Future studies should examine the aforementioned mechanism for turn-taking by conducting an online
789 experiment of human-robot interactions. However, the computational cost of online error-regression for
790 the posterior inference has been the major bottleneck for conducting such experiments in real time, and
791 this is why the current study was limited to a simulation of pseudo-imitative interaction using recorded
792 visuo-proprioceptive sequence patterns, rather than introducing actual, real-time, human-robot interaction.
793 Although our group has shown that some real-time experiments using online ER are possible using only
794 the sensory modality of proprioception (Chame and Tani, 2019), it becomes prohibitive when also using
795 vision, with sufficient pixel resolution. Regarding this problem, some may suggest employing other types
796 of variational models, such as a variational recurrent neural network (VRNN) (Chung et al., 2015), because
797 a VRNN demands far less computation time, since the posterior at each time step can be inferred by simple
798 sequential mapping of inputs using an autoencoder (Kingma et al., 2016). However, the current scheme for
799 inference of the posterior through iterative computation for optimization is probably vital for any embodied
800 cognitive systems that require rapid adaptation of internal states to the environment. Actually, Ahmadi and
801 Tani (Ahmadi and Tani, 2019) showed that PV-RNN performs better than VRNN in online prediction in
802 dynamically changing environments by inferring the posterior using the error-regression scheme. Therefore,
803 future studies should explore possible methods for accelerating online error-regression of the model, such
804 as by massive parallelization so as to conduct real-time, human-robot interactions using the current model.

ETHICS STATEMENT

805 Written informed consent was obtained from the individuals for publication of any potentially identifiable
806 images or data included in this article.

CONFLICT OF INTEREST STATEMENT

807 The authors declare that the research was conducted in the absence of any commercial or financial
808 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

809 WO and JT conceived the concepts and models and contributed to the writing. WO conducted the
810 experiments.

FUNDING

811 This study was supported by funding from Okinawa Institute of Science and Technology Graduate
812 University. This study has also been partially supported by a Grant-in-Aid for Scientific Research(A) in
813 Japan, 20H00001, “Phenomenology of Altered Consciousness: An Interdisciplinary Approach through
814 Philosophy, Mathematics, Neuroscience, and Robotics”.

ACKNOWLEDGMENTS

815 We thank lab members in the Cognitive Neurorobotics Research Unit. We are especially grateful to
816 Ahmadreza Ahmadi and Prasanna Vijayaraghavan for their help in developing the model. We thank Siqing
817 Hou for his help in collecting data. We also thank Steven D. Aird for editing the manuscript.

REFERENCES

- 818 Ahmadi, A. and Tani, J. (2017). How can a recurrent neurodynamic predictive coding model cope with
819 fluctuation in temporal patterns? robotic experiments on imitative interaction. *Neural Networks* 92, 3–16
- 820 Ahmadi, A. and Tani, J. (2019). A novel predictive-coding-inspired variational rnn model for online
821 prediction and recognition. *Neural computation* 31, 2025–2074
- 822 Aly, A. and Tapus, A. (2015). An online fuzzy-based approach for human emotions detection: An overview
823 on the human cognitive model of understanding and generating multimodal actions. In *Intelligent*
824 *assistive robots* (Springer). 185–212
- 825 Baltieri, M. and Buckley, C. L. (2017). An active inference implementation of phototaxis. In *Artificial Life*
826 *Conference Proceedings 14* (MIT Press), 36–43
- 827 Battaglia, P. W., Jacobs, R. A., and Aslin, R. N. (2003). Bayesian integration of visual and auditory signals
828 for spatial localization. *Josa a* 20, 1391–1397
- 829 Boucenna, S., Cohen, D., Meltzoff, A. N., Gaussier, P., and Chetouani, M. (2016). Robots learn to
830 recognize individuals from imitative encounters with people and avatars. *Scientific reports* 6, 19908
- 831 Boucenna, S., Gaussier, P., Andry, P., and Hafemeister, L. (2014). A robot learns the facial expressions
832 recognition and face/non-face discrimination through an imitation game. *International Journal of Social*
833 *Robotics* 6, 633–652
- 834 Buckley, C. L., Kim, C. S., McGregor, S., and Seth, A. K. (2017). The free energy principle for action and
835 perception: A mathematical review. *Journal of Mathematical Psychology* 81, 55–79
- 836 Chame, H. F. and Tani, J. (2019). Cognitive and motor compliance in intentional human-robot interaction.
837 *arXiv preprint arXiv:1911.01753* Accepted for publication in IEEE ICRA2020

- 838 Choi, M. and Tani, J. (2018). Predictive coding for dynamic visual processing: Development of functional
839 hierarchy in a multiple spatiotemporal scales rnn model. *Neural computation* 30, 237–270
- 840 Chung, J., Kastner, K., Dinh, L., Goel, K., Courville, A. C., and Bengio, Y. (2015). A recurrent latent
841 variable model for sequential data. In *Advances in neural information processing systems*. 2980–2988
- 842 Clark, A. (2015). *Surfing uncertainty: Prediction, action, and the embodied mind* (Oxford University
843 Press)
- 844 Di Pellegrino, G., Fadiga, L., Fogassi, L., Gallese, V., and Rizzolatti, G. (1992). Understanding motor
845 events: a neurophysiological study. *Experimental brain research* 91, 176–180
- 846 Friston, K. (2005). A theory of cortical responses. *Philosophical transactions of the Royal Society B:*
847 *Biological sciences* 360, 815–836
- 848 Friston, K. (2012). Prediction, perception and agency. *International Journal of Psychophysiology* 83,
849 248–252
- 850 Friston, K. (2018). Does predictive coding have a future? *Nature neuroscience* 21, 1019
- 851 Friston, K., Mattout, J., and Kilner, J. (2011). Action understanding and active inference. *Biological*
852 *cybernetics* 104, 137–160
- 853 Friston, K. J., Daunizeau, J., and Kiebel, S. J. (2009). Reinforcement learning or active inference? *PloS*
854 *one* 4, e6421
- 855 Friston, K. J., Daunizeau, J., Kilner, J., and Kiebel, S. J. (2010). Action and behavior: a free-energy
856 formulation. *Biological cybernetics* 102, 227–260
- 857 Gallagher, S. (2000). Philosophical conceptions of the self: implications for cognitive science. *Trends in*
858 *cognitive sciences* 4, 14–21
- 859 Gallese, V., Fadiga, L., Fogassi, L., and Rizzolatti, G. (1996). Action recognition in the premotor cortex.
860 *Brain* 119, 593–609
- 861 Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., et al. (2017). beta-vae: Learning
862 basic visual concepts with a constrained variational framework. *Iclr* 2, 6
- 863 Hohwy, J. (2013). *The predictive mind* (Oxford University Press)
- 864 Hurley, S. L. (2005). *Perspectives on imitation: From neuroscience to social science* (MIT press)
- 865 Huys, R., Daffertshofer, A., and Beek, P. J. (2004). Multiple time scales and multiform dynamics in
866 learning to juggle. *Motor control* 8, 188–212
- 867 Hwang, J., Kim, J., Ahmadi, A., Choi, M., and Tani, J. (2020). Dealing with large-scale spatio-temporal
868 patterns in imitative interaction between a robot and a human by using the predictive coding framework.
869 *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 50, 1918–1931
- 870 Ikegami, T. and Iizuka, H. (2007). Turn-taking interaction as a cooperative and co-creative process. *Infant*
871 *Behavior and Development* 30, 278–288
- 872 Ito, M. and Tani, J. (2004). On-line imitative interaction with a humanoid robot using a dynamic neural
873 network model of a mirror system. *Adaptive Behavior* 12, 93–115
- 874 Kawato, M., Furukawa, K., and Suzuki, R. (1987). A hierarchical neural-network model for control and
875 learning of voluntary movement. *Biological cybernetics* 57, 169–185
- 876 Kilner, J. M., Friston, K. J., and Frith, C. D. (2007). Predictive coding: an account of the mirror neuron
877 system. *Cognitive processing* 8, 159–166
- 878 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint*
879 *arXiv:1412.6980*
- 880 Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., and Welling, M. (2016). Improved
881 variational inference with inverse autoregressive flow. In *Advances in neural information processing*
882 *systems*. 4743–4751

- 883 Kohler, E., Keysers, C., Umiltà, M. A., Fogassi, L., Gallese, V., and Rizzolatti, G. (2002). Hearing sounds,
884 understanding actions: action representation in mirror neurons. *Science* 297, 846–848
- 885 Kording, K. P., Tenenbaum, J. B., and Shadmehr, R. (2007). The dynamics of memory as a consequence of
886 optimal adaptation to a changing body. *Nature neuroscience* 10, 779–786
- 887 LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., et al. (1989).
888 Backpropagation applied to handwritten zip code recognition. *Neural computation* 1, 541–551
- 889 LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document
890 recognition. *Proceedings of the IEEE* 86, 2278–2324
- 891 Lee, T. S. and Mumford, D. (2003). Hierarchical bayesian inference in the visual cortex. *JOSA A* 20,
892 1434–1448
- 893 Legaspi, R. and Toyozumi, T. (2019). A bayesian psychophysics model of sense of agency. *Nature*
894 *communications* 10, 1–11
- 895 Liu, R., Lehman, J., Molino, P., Such, F. P., Frank, E., Sergeev, A., et al. (2018). An intriguing failing
896 of convolutional neural networks and the coordconv solution. In *Advances in Neural Information*
897 *Processing Systems*. 9605–9616
- 898 Moore, J. W., Wegner, D. M., and Haggard, P. (2009). Modulating the sense of agency with external cues.
899 *Consciousness and cognition* 18, 1056–1064
- 900 Nadel, J. (2002). Imitation and imitation recognition: Functional use in preverbal infants and nonverbal
901 children with autism. *The imitative mind: Development, evolution, and brain bases* 4262
- 902 Newell, K. M., Liu, Y.-T., and Mayer-Kress, G. (2001). Time scales in motor learning and development.
903 *Psychological review* 108, 57
- 904 Nishimoto, R. and Tani, J. (2009). Development of hierarchical structures for actions and motor imagery: a
905 constructivist view from synthetic neuro-robotics study. *Psychological Research PRPF* 73, 545–558
- 906 Ogata, T., Nishide, S., Kozima, H., Komatani, K., and Okuno, H. G. (2010). Inter-modality mapping in
907 robot with recurrent neural network. *Pattern Recognition Letters* 31, 1560–1569
- 908 Oliver, G., Lanillos, P., and Cheng, G. (2019). Active inference body perception and action for humanoid
909 robots. *arXiv preprint arXiv:1906.03022*
- 910 Oztop, E., Kawato, M., and Arbib, M. (2006). Mirror neurons and imitation: A computationally guided
911 review. *Neural networks* 19, 254–271
- 912 Oztop, E., Kawato, M., and Arbib, M. A. (2013). Mirror neurons: functions, mechanisms and models.
913 *Neuroscience letters* 540, 43–55
- 914 Pezzulo, G., Rigoli, F., and Friston, K. J. (2018). Hierarchical active inference: A theory of motivated
915 control. *Trends in cognitive sciences* 22, 294–306
- 916 Pitti, A., Quoy, M., Lavandier, C., and Boucenna, S. (2020). Gated spiking neural network using iterative
917 free-energy optimization and rank-order coding for structure learning in memory sequences (inferno
918 gate). *Neural Networks* 121, 242–258
- 919 Rao, R. P. and Ballard, D. H. (1999). Predictive coding in the visual cortex: a functional interpretation of
920 some extra-classical receptive-field effects. *Nature neuroscience* 2, 79
- 921 Rizzolatti, G. and Fogassi, L. (2014). The mirror mechanism: recent findings and perspectives.
922 *Philosophical Transactions of the Royal Society B: Biological Sciences* 369, 20130420
- 923 Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985). *Learning internal representations by error*
924 *propagation*. Tech. rep., California Univ San Diego La Jolla Inst for Cognitive Science
- 925 Sacks, H., Schegloff, E. A., and Jefferson, G. (1978). A simplest systematics for the organization of turn
926 taking for conversation. In *Studies in the organization of conversational interaction* (Elsevier). 7–55

- 927 Shimojo, S. (2014). Postdiction: its implications on visual awareness, hindsight, and sense of agency.
928 *Frontiers in psychology* 5, 196
- 929 Smith, M. A., Ghazizadeh, A., and Shadmehr, R. (2006). Interacting adaptive processes with different
930 timescales underlie short-term motor learning. *PLoS biology* 4
- 931 Synofzik, M., Vosgerau, G., and Newen, A. (2008). Beyond the comparator model: a multifactorial
932 two-step account of agency. *Consciousness and cognition* 17, 219–239
- 933 Tani, J. and Nolfi, S. (1999). Learning to perceive the world as articulated: an approach for hierarchical
934 learning in sensory-motor systems. *Neural Networks* 12, 1131–1141
- 935 Trevarthen, C. (1979). Communication and cooperation in early infancy: A description of primary
936 intersubjectivity. *Before speech: The beginning of interpersonal communication* 1, 530–571
- 937 Valentin, P., Boucenna, S., Gaussier, P., and Pitti, A. (2019). Robot recognizing vowels in a multimodal
938 way. In *2019 Joint IEEE 9th International Conference on Development and Learning and Epigenetic
939 Robotics (ICDL-EpiRob)* (Oslo, Norway), 9th International Conference on Development and Learning
940 and Epigenetic Robotics (ICDL-EpiRob), 103–104
- 941 Werbos, P. (1974). Beyond regression: New tools for prediction and analysis in the behavioral sciences.
942 *Ph. D. dissertation, Harvard University*
- 943 Yamashita, Y. and Tani, J. (2008). Emergence of functional hierarchy in a multiple timescale neural
944 network model: a humanoid robot experiment. *PLoS computational biology* 4, e1000220