# A Neurorobotics Simulation of Autistic Behavior Induced by Unusual Sensory Precision

Hayato Idei      Shingo Murata      Yiwen Chen      Yuichi Yamashita      Jun Tani

Tetsuya Ogata*

## Abstract

Recently, applying computational models developed in cognitive science to psychiatric disorders has been recognized as an essential approach for understanding cognitive mechanisms underlying psychiatric symptoms. Autism spectrum disorder is a neurodevelopmental disorder that is hypothesized to affect information processes in the brain involving the estimation of sensory precision (uncertainty), but the mechanism by which observed symptoms are generated from such abnormalities has not been thoroughly investigated. Using a humanoid robot controlled by a neural network using a precision-weighted prediction error minimization mechanism, it is suggested that both increased and decreased sensory precision could induce the behavioral rigidity characterized by resistance to change that is characteristic of autistic behavior. Specifically, decreased sensory precision caused any error signals to be disregarded, leading to invariability of the robot's intention, while increased sensory precision caused an excessive response to error signals, leading to fluctuations and subsequent fixation of intention. The results may provide a system-level explanation of mechanisms underlying

different types of behavioral rigidity in autism spectrum and other psychiatric disorders. In addition, our findings suggest that symptoms caused by decreased and increased sensory precision could be distinguishable by examining the internal experience of patients and neural activity coding prediction error signals in the biological brain.

*Correspondence concerning this article should be addressed to Tetsuya Ogata, Department of Intermedia Art and Science, Waseda University, Tokyo, Japan. E-mail: ogata@waseda.jp. H. Idei, Department of Intermedia Art and Science, Waseda University, Tokyo, Japan; S. Murata, Y. Chen, Department of Modern Mechanical Engineering, Waseda University, Tokyo, Japan; Y. Yamashita, Department of Functional Brain Research, National Center of Neurology and Psychiatry, Tokyo, Japan; J. Tani, Cognitive Neurorobotics Research Unit, Okinawa Institute of Science and Technology (OIST), Okinawa, Japan.

## 1   Introduction

Autism spectrum disorder (ASD) is a neurodevelopmental disorder that affects a broad range of cognitive functions, including perception[1], action[2], and social cognition[3]. In particular, behavioral rigidity manifested as restricted, repetitive behavior and resistance to change is a core ASD symptom[4–7], albeit such behavioral rigidity can be also observed in other psychiatric disorders[8,9]. Behavioral rigidity in ASD consists of various behavioral categories, such as stereotyped motor mannerisms (e.g., hand flapping) and self-injurious or compulsive behavior[10,11]. Although the reduced behavioral flexibility severely limits the social adaptation of patients, its cause and the underlying cognitive mechanisms remain unclear.

There have been many studies aiming to construct theories that explain the mechanisms underlying autistic symptoms[3,12,13], and recently the focus of these attempts has shifted to the idea of describing fundamental brain function as a set of computational processes[14]. In particular, theoretical explanations based on prediction error minimization frameworks, such as predictive coding[15,16] and the free energy principle[17], have been well investigated because they may be able to uniformly explain various ranges of

autistic symptoms using a simple and neurologically plausible principle[18–24]. The prediction error minimization mechanism explains how we acquire knowledge and skills (learning), and how we successively infer the causes of sensory inputs and recognize environments as the process of updating a model of the world based on minimizing error between a prediction about incoming sensory inputs and actual sensory inputs. Within a scheme in which prediction error causes the brain to update its model of the world, it is crucial to estimate precision (inverse variance) of sensory information: the expected precision of certain sensory information can provide information about the reliability of the generated prediction error, which influences how much weight is given to the error when updating predictions. For example, although prediction errors for certain sensory inputs that contain information refuting the current expectation (e.g., one looks around the seabed in clear water and what seems like sand suddenly moves) should cause brain to update its expectation (one recognizes it is not sand but flatfish), errors in sensory inputs that are very noisy (one looks around the seabed sand of foggy water and something moves) should not cause the update (one would think it is only a wave causing the movement). Although the estimation of such context-dependent sensory precision (prediction about whether information is informative or just noise) helps us to be flexible and adaptable in an uncertain world, deficits of it are expected to cause perceptual peculiarity and great difficulty in social contexts that are filled with situations of particularly high complexity and uncertainty[21–25]. Van de Cruys et al.[21] suggested that inflexibly overestimated sensory precision causes autistic symptoms and inflexible behavior may be considered as an attempt to minimize prediction errors; otherwise, patients are exposed to huge error signals. Lawson et al.[22] explains autistic behaviors as the consequences of "an imbalance of the precision ascribed to sensory evidence relative to prior beliefs". These aberrant precision accounts for ASD in previous studies are normative and testable, but only suggestive. Specifically, there is a gap between the cognitive mechanisms described in the theories and the actual generation of the symptoms.

This kind of problem is broadly described in psychiatry, and there is a need to demonstrate actual generation of symptoms using formal computational models[26–30]. Indeed, several computational simulations of psychiatric symptoms have been conducted to try to understand the processes underlying these symptoms and clarify the relationships between abnormalities at neurological- and behavioral levels[31–38]. In particular, embodiment[39,40] in a robot agent acting in physical environments may be useful or even essential for understanding the cognitive mechanisms of psychiatric disorders. That is because psychiatric disorders are characterized by behavioral and perceptual conditions observed through interaction with real environments and physical agents. In a related study, Yamashita and Tani[34] performed a neurorobotics experiment to investigate schizophrenic cognition by utilizing a hierarchical neural network model. Their robotic experiment showed that behaviors analogous to psychiatric symptoms, such as fictive sensations and cataleptic, stereotyped behaviors, can be generated in the coupled dynamics describing the neural networks, body and environment due to synaptic disconnections between different levels of the neural network.

In this study we investigated the effects of increased and decreased sensory precision on adaptive behaviors by conducting experiments using a humanoid robot implemented with a version of the predictive coding model. In the experiment, a task involving adaptive interaction between the robot and a human experimenter was considered. Initially, the neural network model inside the robot learned to generate a set of sequence patterns representing different behaviors of the robot. After the learning phase, the level of estimated sensory precision was manipulated. Then, the change in the robot's behavior in response to the alteration of the level of sensory precision was observed through experiments in which the robot was required to appropriately recognize situations determined by the experimenter. The results show both increased and decreased sensory precision can cause seemingly similar inflexible behavioral patterns, such as inappropriate repetitive behavior and freezing; but these behaviors are the result of different processes at the network-level in the two cases. Our findings may

provide a system-level account for different types of behavioral rigidity observed in ASD and other psychiatric disorders, and extends computational perspectives on the cognitive mechanisms underlying psychiatric symptoms.

## 2 Methods

### 2.1 Computational framework

We used an artificial recurrent neural network (RNN) model to investigate the effects of increased and decreased sensory precision on adaptive behaviors of a robot. An RNN is a connectionist model which can process temporal sequences thanks to recurrent connections between neural units[41]. Owing to its capacity to learn to reproduce complex dynamic behaviors, RNNs have been used in cognitive neurorobotics studies aiming to understand human cognition[42, 43]. Murata et al.[44], within the cognitive robotics scheme, proposed an RNN model with a mechanism for estimating the time-varying uncertainty of sensory information in terms of variance (inverse precision) as inspired by the free energy minimization principle proposed by Friston[17]. This RNN, called a stochastic-continuous time RNN (S-CTRNN), can learn to predict not only sensory inputs but also their variances based on negative log-likelihood minimization, which is equivalent to precision-weighted prediction error minimization. Tani[45] proposed an RNN with parametric bias (RNNPB) which has an online adaptation mechanism based on prediction error minimization. In this framework, PB is encoded in a small group of neural units which works as a higher-level neural representation of the network behavior, and the associations between specific patterns of PB activity and different temporal training patterns are self-organized through a learning process. Owing to this characteristic of PB, a robot driven by RNNPB can not only generate multiple learned behavioral patterns but also switch its behavior by adaptively modulating the PB states in response to a discrepancy between a prediction and actual sensory information. PB states thus can be regarded as the higher-level "intention" of a robot. Utilizing this model, Ito et al.[46] demonstrated flexible switching of ball-playing behaviors by a humanoid robot in response to changes in the environment.

In the present study, an S-CTRNN with PB was adopted as the computational model for simulating aberrant sensory precision because of its capacity to learn to estimate sensory variance (precision) and adapt to different environments using a prediction error minimization mechanism. The following subsections describe in detail the mathematical procedures used for the forward dynamics and parameter optimization of the S-CTRNN with PB.

#### 2.1.1 Forward dynamics

The neuronal model is a conventional firing rate model. The internal state of the $i$th neural unit at time step $t$, $u_{t,i}^{(s)}$ $(t \geq 1)$, is described by

$$
u_{t,i}^{(s)} = \begin{cases} u_{t-1,i}^{(s)} & (i \in I_{\mathrm{P}}), \\[2ex] \dfrac{1}{\tau_i}\left(\displaystyle\sum_{j \in I_{\mathrm{I}}} w_{ij} x_{t,j}^{(s)} + \sum_{j \in I_{\mathrm{C}}} w_{ij} c_{t-1,j}^{(s)} + \sum_{j \in I_{\mathrm{P}}} w_{ij} p_{t,j}^{(s)} \right. \\[2ex] \left. + b_i \right) + \left(1 - \dfrac{1}{\tau_i}\right) u_{t-1,i}^{(s)} & (i \in I_{\mathrm{C}}), \\[2ex] \displaystyle\sum_{j \in I_{\mathrm{C}}} w_{ij} c_{t,j}^{(s)} + b_i & (i \in I_{\mathrm{O}}, I_{\mathrm{V}}). \end{cases}
$$

$$(1)$$

Here, $I_{\mathrm{I}}$, $I_{\mathrm{P}}$, $I_{\mathrm{C}}$, $I_{\mathrm{O}}$, and $I_{\mathrm{V}}$ are index sets of the input, PB, context, output, and variance neural units, respectively, $w_{ij}$ is the weight of the synaptic connection from the $j$th neuron to the $i$th neuron, $x_{t,j}^{(s)}$ is the $j$th input of the $s$th sequence at time step $t$, $c_{t,j}^{(s)}$ is the $j$th context state, $p_{t,j}^{(s)}$ is the $j$th PB state, $b_i$ is the bias of the $i$th neuron, and $\tau_i$ is the time constant of the $i$th neuron. From this equation, we see that PB units can be considered to be a specific type of context unit whose time constant is infinite.
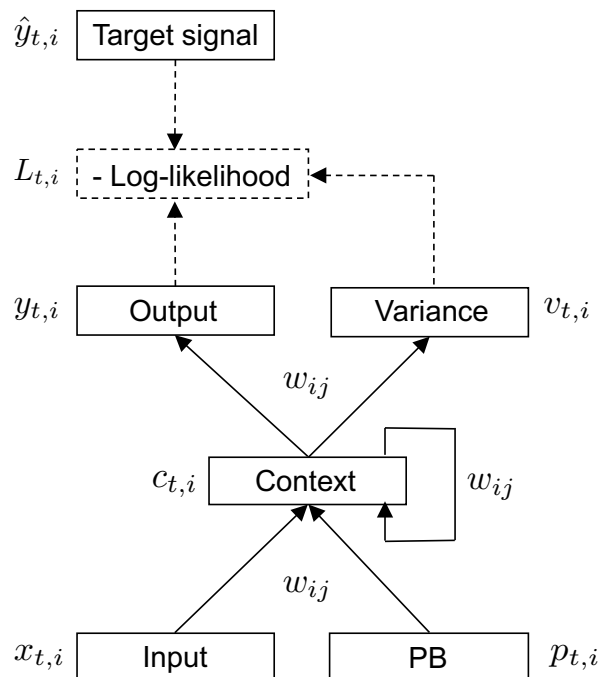
$$\hat{y}_{t,i} \quad \boxed{\text{Target signal}}$$

$$L_{t,i} \quad \boxed{\text{- Log-likelihood}}$$

$$y_{t,i} \quad \boxed{\text{Output}} \qquad \boxed{\text{Variance}} \quad v_{t,i}$$

$$w_{ij}$$

$$c_{t,i} \quad \boxed{\text{Context}} \quad w_{ij}$$

$$w_{ij}$$

$$x_{t,i} \quad \boxed{\text{Input}} \qquad \boxed{\text{PB}} \quad p_{t,i}$$

Figure 1: The S-CTRNN utilized in this study has five groups of neural units: input, context, output, variance, and PB units. Input neural units receive current sensory inputs $x_t$. Based on the inputs, PB state $p_t$, and context state $c_t$, the S-CTRNN generates predictions about the mean $y_t$ and variance $v_t$ of future inputs in the output and variance units, respectively. Parameters, such as synaptic weights $w_{ij}$ and the internal state of PB units, are optimized by minimizing negative log-likelihood as calculated using predictions about sensory states, their variance, and actual target sensory states $\hat{y}_t$.

The current study sets all initial values of the internal states of the context units to zero, while those of the PB units are optimized for each target sequence in the learning phase. This indicates that differences between target sequences are represented in the activity of the PB units.

The activation values of each neural unit are calculated as

$$p_{t,i}^{(s)} = \tanh\left(u_{t,i}^{(s)}\right) \qquad (0 \le t \wedge i \in I_\mathrm{P})\,, \quad (2)$$

$$c_{t,i}^{(s)} = \tanh\left(u_{t,i}^{(s)}\right) \qquad (0 \le t \wedge i \in I_\mathrm{C})\,, \quad (3)$$

$$y_{t,i}^{(s)} = \tanh\left(u_{t,i}^{(s)}\right) \qquad (1 \le t \wedge i \in I_\mathrm{O})\,, \quad (4)$$

$$v_{t,i}^{(s)} = \exp\left(u_{t,i}^{(s)}\right) \qquad (1 \le t \wedge i \in I_\mathrm{V})\,. \quad (5)$$

### 2.1.2 Parameter optimization

The neural network performs parameter optimization based on the gradient decent method aiming to minimize the objective function

$$L_{t,i}^{(s)} = \frac{\ln\left(2\pi v_{t,i}^{(s)}\right)}{2} + \frac{\left(\hat{y}_{t,i}^{(s)} - y_{t,i}^{(s)}\right)^2}{2v_{t,i}^{(s)}}, \qquad (6)$$

where $\hat{y}_{t,i}^{(s)}$ is the $i$th target value corresponding to the $s$th sequence. Minimizing this negative log-likelihood can be regarded as minimizing the precision-weighted (inverse variance-weighted) prediction error, and is formally equivalent to minimizing free energy in the active inference scheme proposed by Friston[17].

In the learning phase, parameters, including synaptic weights $w_{ij}$, biases $b_i$, and the initial internal

states of PB units $u_{0,i}^{(s)}$ $(i \in I_{\mathrm{P}})$, are updated in an offline manner. Parameter optimization is performed by minimizing the sum of the negative log-likelihood over all dimensions, time steps, and sequences as

$$L = \sum_{s \in I_{\mathrm{S}}} \sum_{t=1}^{T^{(s)}} \sum_{i \in I_{\mathrm{O}}} L_{t,i}^{(s)}, \qquad (7)$$

where $I_{\mathrm{S}}$ and $T^{(s)}$, respectively, represent the index set and the length of the $s$th target sequence. The partial derivative of each parameter, $(\partial L / \partial \boldsymbol{\theta})$, can be found using the back-propagation-through-time (BPTT) method described in[44, 47].

In the adaptation phase, after learning, only the internal states of the PB units are optimized online and other parameters are fixed. The negative log-likelihood within a short time window $W$ is accumulated as

$$L = \sum_{t'=t-W+1}^{t} \sum_{i \in I_{\mathrm{O}}} L_{t',i}^{(s)}. \qquad (8)$$

The time window of length $W$ moves along with the increment of the network time step $t$. Using the accumulated negative log-likelihood, the internal states of the PB units at time step $t - W$ are optimized. The partial derivatives of the internal states of PB units are also calculated by the BPTT algorithm.

In both the learning and adaptation phases, parameters which are allowed to be optimized are collected as a vector $\boldsymbol{\theta}$, and $\boldsymbol{\theta}$ at the $n$th epoch is updated using gradient decent on the accumulated negative log-likelihood $L$:

$$\boldsymbol{\theta}(n) = \boldsymbol{\theta}(n-1) + \Delta\boldsymbol{\theta}(n), \qquad (9)$$

$$\Delta\boldsymbol{\theta}(n) = -\alpha \frac{\partial L}{\partial \boldsymbol{\theta}} + \eta \Delta\boldsymbol{\theta}(n-1). \qquad (10)$$

Here, $\alpha$ is the learning rate, and $\eta$ is a coefficient representing the momentum term. In this study, $\alpha$ and $\eta$ are set at 0.0001 and 0.9, respectively.

## 2.2 Task setting

In order to provide the robot with a task suitable for testing our hypothesis that aberrant sensory precision induces behavioral rigidity, we require a dynamical interaction setup in which the robot needs to perceive sensory information with intrinsic uncertainty and flexibly recognize situations determined by others. We chose a ball-playing scheme involving interaction between a robot and a human experimenter that was used in a previous study by Chen et al.[48]. The behavioral patterns of the robot consist of four different ball-playing behaviors (See Fig. 2a). In the "right" and "left" behaviors, the robot is required to wait for the ball coming from the human subject and then return it. "Self-play" behavior consists of rolling the ball in front of itself, and the "attract" behavior is an up-down motor action with the arms while the partner engages in the "self-play" behavior of moving the ball left and right. After the S-CTRNN with PB learned to reproduce these visuo-proprioceptive temporal patterns, the behavioral performance of the robot with the trained neural network model was tested in the task of adaptive ball-playing interaction with a human subject.

## 2.3 Experimental environment

We employed a small humanoid robot NAO (Aldebaran) that has a body corresponding to only the upper half of the human body. The robot sat in front of a workbench and engaged in a ball-playing interaction with a human experimenter standing on the opposite side of the bench. The robot's action involved only movements of the arms with 4 degrees of freedom for each arm (2 shoulders and 2 elbows). In addition, a camera located in the robot's mouth obtained the center of gravity coordinates for the yellow object which was used as 2-dimensional inputs for ball position. Using the minimum and maximum values of each piece of sensory information, the values of joint angles and the ball position were mapped to values ranging from -0.8 to 0.8. The size of the workbench and the diameter of the ball are approximately $45 \times 5 \times 30$ cm and 9 cm, respectively.

## 2.4 Training

Training of the neural network was conducted in an offline manner by supervised learning using target perceptual sequences recorded in advance. The target perceptual sequences were recorded while the
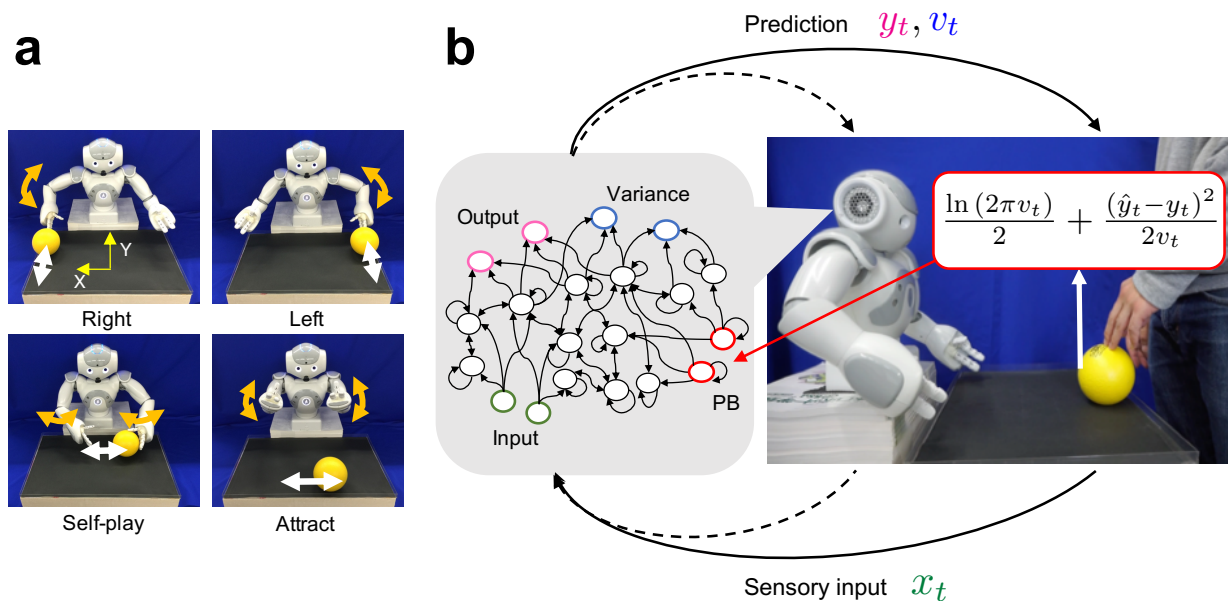
Figure 2: (a) Four interactive behavioral patterns learned by a robot controlled by an S-CTRNN with PB. The upper-left and -right figures show the right and left behaviors, respectively. The lower-left and -right figures show the self-play and attract behaviors. (b) System overview during adaptive interaction between a robot and a experimenter. The solid lines for prediction and sensory input represent visual information about the ball position. The dotted lines represent proprioceptive information about the robot's joint angles. The neural network generates predictions about sensory states $y_t$ and their variances $v_t$ based on current sensory inputs $x_t$, and also recognizes situations by updating PB activity online in the direction of minimizing the negative log-likelihood calculated using the predictions and the target signal (actual sensory feedback) $\hat{y}_t$.

robot repeatedly performed each ball-playing behavior, where the arm movement was generated exactly following preprogrammed trajectories instead of the ones generated by the neural network model. Each of the 4 behavioral patterns was obtained as a sequence of 10-dimensional vectors (8 dimensions for joint angles and 2 dimensions for ball position). For the training, 3 sequences were prepared for each behavioral pattern. The time lengths of the sequences were approximately 1600 time steps for "right", 1900 time steps for "left", 1600 time steps for "self-play", and 1200 time steps for "attract".

The neural network learned to reproduce these target visuo-proprioceptive sequences. The objective of the learning is to find the optimal values of the parameters (synaptic weights, biases, and internal states of PB units) minimizing negative log-likelihood, or precision-weighted prediction error. At

first, each parameter was initialized with a random value and the network produced random sequences. The parameters were updated in the direction of minimizing negative log-likelihood accumulated through the duration of the target sequences. Repeating the update process many times, the network became able to produce visuo-proprioceptive sequences with the same stochastic properties as the target sequences. In addition, the associations between a particular pattern of target sequence and specific internal states of PB units self-organized.

## 2.5   Online adaptation

After the learning process, the robot engaged in an adaptive interaction with a human experimenter by updating PB states (intention) online. In this phase, the robot's intention was first set to a certain state

corresponding to a learned behavior, and situation (ball dynamics pattern) was controlled by the experimenter. The goal of the robot was to flexibly recognize situations using visual cues. Real-time adaptation during task execution by the robot was performed based on an interaction between a top-down prediction generation process and a bottom-up parameter adaptation process. In the top-down prediction generation process, the network generated a temporal sequence corresponding to time steps from $t - W + 1$ to $t$, based on the sensory inputs at time step $t - W + 1$ and the constant PB states (intention). The visuo-proprioceptive sequence was generated by a "closed-loop" procedure, meaning that predictions about mean values of the sensory states at a certain time step were used as inputs at the next step. The initial inputs for proprioceptive states at time step $t - W + 1$ were taken from the generated mean predictions at $t - W$, and those for vision states were taken from the vision data caught by the camera at time step $t - W + 1$. In the bottom-up adaptation process, the negative log-likelihood at each time step within time window $W$ was calculated by using the predictions about vision states, their variance, and the actual visual feedback (see Fig. 2b). The PB states (intention) were updated in the direction of minimizing the accumulated negative log-likelihood. Based on the updated PB states, the temporal sequence within the time window was re-generated. After repeating these top-down and bottom-up processes for a certain number of times, the network generated its predictions for time step $t + 1$ and the predictions about proprioceptive states are sent to the robot as the target for subsequent joint positions. This procedure, where recognition and prediction in the past are reconstructed based on current sensory information, is more properly regarded as a "postdiction" process[49,50], and generated predictions for time steps from $t - W + 1$ to $t$ are more suitably referred to as postdiction of the past rather than prediction in the literal sense.

## 2.6 Parameter setting for the experiment

The number of input, output, and variance neural units were $N_I = N_O = N_V = 10$, corresponding to the dimension of the robot's sensory states, and the number of PB units was $N_P = 2$. The number and time constant of the context units were $N_C = 50$ and $\tau_i = 4$, respectively. In the learning phase, the weights of synaptic connections $w_{ij}$ ($j \in I_I, I_C$) and biases $b_i$ were initialized with random values following uniform distributions on the intervals $[-\frac{1}{N_I}, \frac{1}{N_I}]$ ($j \in I_I$) and $[-\frac{1}{N_C}, \frac{1}{N_C}]$ ($j \in I_C$) for weights, and $[-1, 1]$ for biases, and the internal states of PB units were initialized as 0. These parameters are updated offline $300,000$ times in the learning phase. In the adaptation phase, the internal states of PB units were updated online 20 times, and the length of the time window was $W = 10$.

## 2.7 Simulating aberrant sensory precision

This study simulated increased and decreased sensory precision by altering estimated sensory variance (inverse precision). After the network learned to reproduce the set of behavioral patterns, the activation values of the variance units were modified as

$$ v_{t,i}^{(s)} = \exp\left(u_{t,i}^{(s)} + K\right) + \epsilon \quad (i \in I_V), \qquad (11) $$

where $K$ is a constant determining the level of the estimated variance and $\epsilon$ is its minimum value, set as $0.00001$. $K$ is set as 0 in the normal condition, while $K$ is set to negative values in the decreased sensory variance conditions and positive values in the increased sensory variance conditions ($K \in \{-8, -4, 0, 4, 8\}$).

## 2.8 Analysis of robot's behavior

To judge whether the robot's behavior generated during the test phase is appropriate, the generated time series of joint angles was compared with the target (learned) time series. A simple way to compare

two time series is to calculate the distance between the value at each corresponding pair of time steps within a certain time window. However, this method is not necessarily appropriate for comparing a general characteristic of time series because a phase shift will increase the distance between the series. Here, this would increase the distance even when the robot generates the appropriate action. Thus, this study considered histograms of time-series values within a specified time window and then compared the histogram of the time series generated through the test experiment with the target time series. Because a histogram of time series values can be considered as a probability distribution, two time series can be compared by calculating the Kullback–Leibler (KL) divergence. Although the probability distribution lacks some information regarding temporal ordering, this comparative approach is suitable for our purpose because a general characteristic of a time series can be extracted. By considering the amount of the state change and calculating the KL divergence from the learned time series, the behaviors observed in the experiments could be classified into one of four types: "outwardly normal", "freezing" (maintaining one posture), "unlearned movement" (engaging in an unlearned action), and "inappropriate learned movement" (engaging in a learned action other than the target action). These are explained in more detail in below.

To assess the robot's behavior in the experiment, an 8-dimensional time series of joint angles was reduced to a 2-dimensional time series by applying principal component analysis. To extract the probability distribution of the 2-dimensional time series, the 2-dimensional space $[-N, N] \times [-N, N]$ (with $N$ the maximum of the absolute value of time series $S(t) = \{z_1(t), z_2(t)\}$ across all data, where $z_1$ and $z_2$ represent the first and second principal components, respectively) is divided into $N_{\text{bin}}^2$ subspaces (here, $N_{\text{bin}} = 20$). Then, the occurrence frequencies of states within the time series were counted. Based on the acquired probability distributions of the time series, the KL divergence between the probability distribution of the time series generated in the test experiment and the target (learned) time series was calculated. The robot's behavior is judged as

"outwardly normal" if the KL divergence is less than a threshold $\xi$, set here as half of the minimum of KL divergence between each pair of learned time series.

$$D_{\text{KL}}(p\|q) < \xi = \\ 0.5 \times \min_{q_i, q_j \in U_{\hat{s}} \wedge q_i \neq q_j} D_{\text{KL}}(q_i\|q_j). \tag{12}$$

Here, $p$ is the probability distribution of the generated time series through the test experiment, $q$ is the probability distribution of the target movement, and $U_{\hat{s}}$ is a set of the probability distributions of each learned movement.

Atypical behaviors can be classified into one of three types of behaviors according to whether the movements were almost stopped and whether they were close to a learned movement other than the target. We call these "freezing" (if $d < 0.02$ and $\forall q \in U_{\hat{s}}, D_{\text{KL}}(p\|q) \geqq \xi$), "unlearned movement" (if $d \geqq 0.02$ and $\forall q \in U_{\hat{s}}, D_{\text{KL}}(p\|q) \geqq \xi$), and "inappropriate learned movement" (if $d \geqq 0.02$ and $\exists q \in U_{\hat{s}}, D_{\text{KL}}(p\|q) < \xi$). In these, $d$ is the amount of the state change defined as

$$d = \frac{1}{T} \sum_{t=0}^{T} \sum_{i \in I_{\text{O}_{\text{joint}}}} |y_{i,t+1} - y_{i,t}|. \tag{13}$$

Here, $T$ is the length of the time series, $I_{\text{O}_{\text{joint}}}$ is the index set of the joint outputs, and $y_{i,t}$ is output of the $i$th output neural unit at time step $t$.

# 3 Results

## 3.1 Open-ended ball interaction

First, we observed the effects of increased or decreased sensory variance (inverse precision) on the robot's behavior through an open-ended ball interaction where situations (ball dynamics patterns) were changed unpredictably by the experimenter. To assess the robot's behaviors, the joint-angle output of the time series was quantitatively assessed every 100 time steps and classified into one of the four types of movements (see Section 2.8). Fig. 3 and Fig. 4 show some representative examples of the robot's behaviors under each condition. Fig. 5 focuses on the
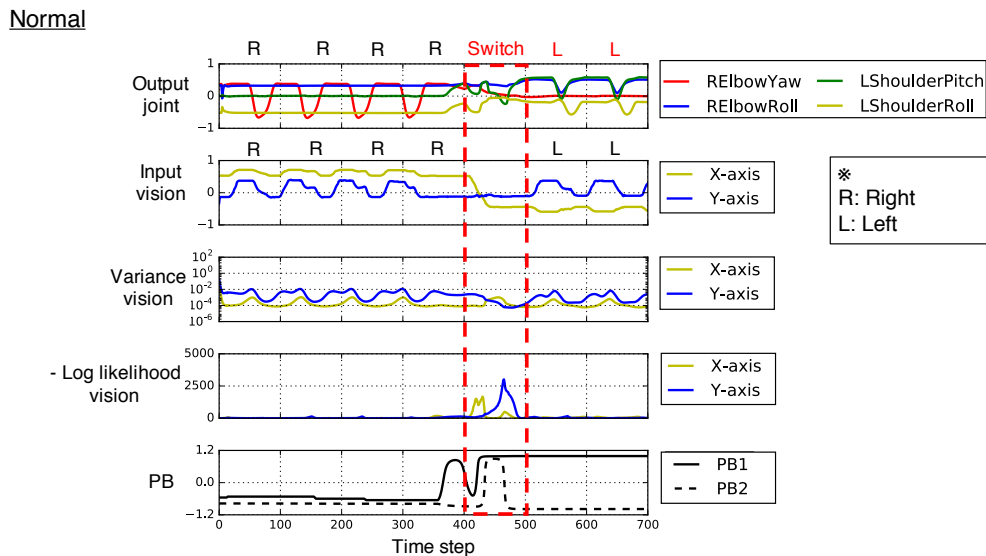
**Figure 3:** Generated time series data from interacting with the experimenter under normal conditions. The robot with a normal network ($K = 0$) successfully adapted to the changing situations (time steps 400–499, in red box) by flexibly switching its intention (PB state) in the direction of minimizing the increased negative log-likelihood. "Output joint" indicates predictions about selected 4-dimensional joint angles. "Input vision", "variance vision", and "negative log-likelihood vision", respectively, indicate the 2-dimensional ball position and corresponding estimated variance and precision-weighted prediction error. The negative log-likelihood at time step $t$ is the value after the postdiction process inside the error regression window between time steps $t - W + 1$ and $t$. PB indicates activation values of the two PB units. The joint-angle output of the time series was quantitatively assessed every 100 time steps as described in Section 2.8.

network-level processes during the trials shown in Fig. 3 and Fig. 4.

Fig. 3 shows a successful interaction between the experimenter and the robot with a normal network. The robot and the experimenter first performed a "right" interaction during time steps 0–399, then the experimenter externally changed the situation (ball dynamics pattern) to a "left" interaction during time steps 400–499 (red box in Fig. 3). The unpredictable situation switch caused conflict between the robot's intention (PB states) and the actual situation. However, the robot's intention was soon updated in the direction of minimizing the increased negative log-likelihood (precision-weighted prediction error) (see also Fig. 5a), and the robot generated behavior appropriate to the situation. This indicates that the robot with a normal network could flexibly recognize and adapt to changing environments.

On the other hand, we observed similar patterns of abnormal overt behaviors, such as freezing or inappropriate repetitive behavior by the robot, under conditions of both increased and decreased sensory variance. Fig. 4a shows freezing behavior under the increased sensory variance condition. In this case, the robot first successfully performed a "right" interaction (time steps 0–399), but the robot almost stopped and maintained a single posture after the situation was switched to "attract" (time steps 500–699). Fig. 4b shows an unlearned repetitive behavior under the decreased sensory variance condition. The robot's action in this case was initially unstable (time steps 0–199) and then converged to an unlearned periodic movement (time steps 200–399), but the robot generated the appropriate movement after the situation was changed (time steps 500–699). These abnormal behaviors, such as freezing and inappropriate repetitive behavior, were observed in both the increased and decreased sensory variance conditions.
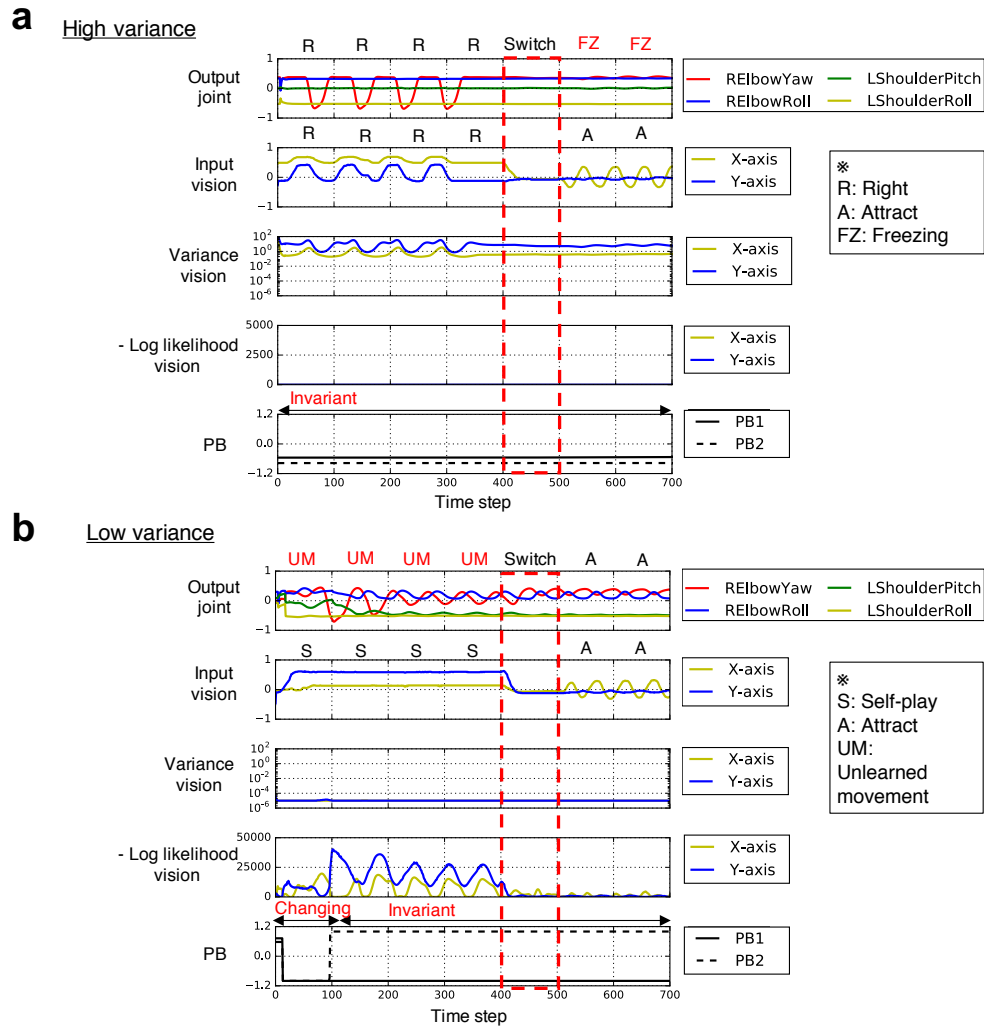
Figure 4: Generated time series data from interacting with the experimenter under increased or decreased sensory variance conditions. (a) Robot's behavior under increased sensory variance condition ($K = 8$). With increased sensory variance, the robot's intention was invariant through the interaction with a situation change (time steps 400–499, in red box) due to highly reduced precision-weighted prediction error, leading to a freezing behavior. (b) Robot's behavior under decreased sensory variance condition ($K = -8$). With decreased sensory variance, the robot experienced huge precision-weighted prediction error signals, and its intention first quickly changed and then fixed at a certain point, leading to an unlearned repetitive movement. Note that the ranges for negative log-likelihood shown in the graphs for the high variance condition and the low variance condition are different. The joint-angle output of the time series was quantitatively assessed every 100 time steps as described in Section 2.8. Abnormal behavioral patterns, including freezing and inappropriate repetitive, were observed under both increased and decreased sensory variance conditions, and these figures show representative examples.
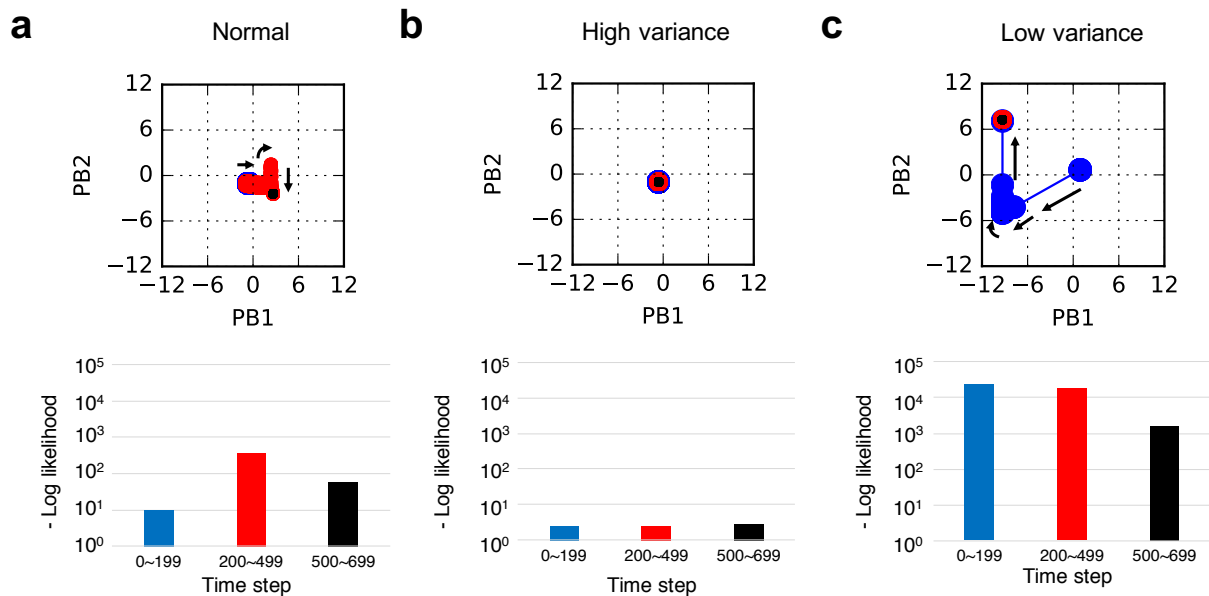
Figure 5: Dynamics of internal PB states (upper figures) and error signals (bottom bar graphs) for each condition during the interactions shown in Fig. 3 and Fig. 4. Colored dots in the upper figures represent PB dynamics during different periods of time (early: time steps 0–199; middle: time steps 200–499; late: time steps 500–699). Bottom bar graphs show the corresponding mean of the negative log-likelihood per time step during each time span. (a) Flexible intention switching under normal condition during the interaction shown in Fig. 3. During the situation change in the middle period, generated error signals caused intention switching and error signals were successfully reduced during interaction in the new situation in the late period. (b) Deficits in intention switching for high sensory variance during the interaction shown in Fig. 4a. Even when the situation changed in the middle period, PB states were almost unchanged due to the under-estimated precision of prediction error. (c) Large shift of network behavior for low sensory variance during the interaction shown in Fig. 4b. Internal PB states first dynamically fluctuated in the early period, but, after the middle period, they became almost fixed at a certain value although generated error signals were still very large.

The videos of the ball interactions and graphs for abnormal behaviors under increased and decreased sensory variance conditions are attached as supplementary information.

In order to distinguish between the mechanisms underlying the similar abnormal behaviors observed in the increased and decreased sensory variance conditions, an analysis was performed on the network-level processes and the level of precision-weighted prediction error the robot experienced (see Fig. 5b and c). In Fig. 5b (see also Fig. 4a), increased sensory variance caused highly reduced precision-weighted prediction error and consequent invariability of the robot's intention (PB states), regardless of the situation change during time steps 400–499 (red box in Fig. 4a). This caused a mismatch between the robot's intention and the situation, leading to freezing behavior. In Fig. 5c (see also Fig. 4b), which shows a decreased sensory variance condition, the internal PB states first quickly but incorrectly changed, possibly because the robot experienced huge precision-weighted prediction errors, which may have included errors associated with inherent noise of the ball dynamics. However, the speed of the changes slowed down when the absolute values of the internal PB states became large. After the repetitive quick state changes and a subsequent slowing down, the internal PB states were fixed at inappropriate values even though the robot was still exposed to error signals as large as, or even larger than, it experienced before the

intentional states became fixed. The fixation of intention caused a mismatch between the robot's intention and the situation, leading to unlearned repetitive behavior. The fixation of PB states may be considered to be the result of fixing at a suboptimal local solution (suboptimal critical point) of the prediction error minimization.

The abnormal behavioral patterns characterized by resistance to change, such as freezing and inappropriate repetitive behavior, may have appeared as a result of the network dynamics converging to fixed points when there was a discrepancy between the robot's intention and the actual situation. In addition to the behavioral abnormalities, generating appropriate behavior in a restricted situation (time steps 0–399 in Fig. 4a and 500–699 in Fig. 4b) was a remarkable characteristic of the observed inflexible behaviors induced by aberrant sensory variance. Thus, the difficulties of the robot should not be attributed to deficits in generating organized behaviors per se, but to deficits in adaptability. This behavioral rigidity characterized by resistance to change may be considered to be analogous to the characteristics of autistic behavior.

## 3.2 Evaluation of adaptability and error signal level

To quantitatively evaluate the frequencies of abnormal overt behaviors described in the previous section, an additional simpler experiment was conducted. In this experiment, the situation set by the experimenter was not changed, but there was a discrepancy between the robot's initial intention (PB states) and the situation. For example, intention of the robot was first set to the value for "left" behavior, but the experimenter rolled the ball to the right. To flexibly interact with the experimenter, the robot thus needed to switch its intention using the visual cue and generate the appropriate behavior. There were six combinations of initial PB states and ball dynamics: initial PB states were "left" or "right" and the experimenter used one of the three other patterns of ball dynamics. Two trials were performed for each combination.

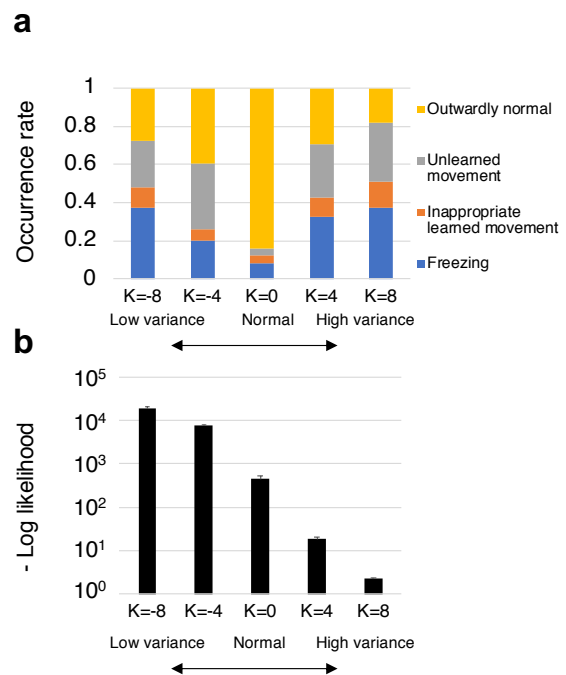We evaluated the robot's behavior in the five con-



Figure 6: Changes in the robot's behavior and negative log-likelihood associated with various levels of sensory variance. (a) The occurrence rates of each behavioral trait over 120 trials for each variance level determined by a parameter $K$ are shown. Behavioral traits observed at time step from 150 to 250 were assessed (see Section 2.8). (b) Negative log-likelihood per time step for each level of sensory variance are shown. Bars in the graph correspond to mean values over 120 trials for each parameter $K$. One-way repeated-measures ANOVAs indicated significant differences between the five conditions for the frequencies of the sum of the three abnormal behaviors ($F(4, 36) = 51.0$, $p < 0.05$) and levels of negative log-likelihood ($F(4, 36) = 110.24$, $p < 0.05$). Adjusting for multiple comparisons using the Holm-Bonferroni method, significant differences were found between the normal condition ($K = 0$) and other unusual variance conditions ($K = -8, -4, 4, 8$) in frequencies of abnormal behaviors (all $p < 0.05$). In addition, significant differences in levels of negative log-likelihood between all pairs were reported (all $p < 0.05$).

ditions ($K = -8, -4, 0, 4, 8$) for ten networks trained with differently randomized initial synaptic weights. Fig. 6 shows the changes in robot's behavior and negative log-likelihood (precision-weighted prediction error) per time step associated with the levels of sensory variance. Behavioral traits observed during

time steps 150–250 were assessed and divided into four overt behavioral patterns ("outwardly normal", "freezing", "unlearned movement", and "inappropriate learned movement"), as described in Section 2.8. Outwardly normal behavior basically means that the robot successfully switched its intention and generated appropriate behavior. However, it also includes behaviors for which the robot's intention was fixed in an inappropriate state due to altered sensory variance but the robot nevertheless managed to generate appropriate behavior using only lower-level network processes based on sensory inputs.

One-way repeated-measures ANOVAs and post hoc multiple comparison adjustments using the Holm-Bonferroni method[51] were conducted for the frequencies of the abnormal overt behaviors and the levels of negative log-likelihood. The level of statistical significance was set at $p < 0.05$. The repeated-measures ANOVAs indicated significant differences among the five conditions in the frequencies of the sum of the three abnormal movements ($F(4, 36) = 51.0$, $p < 0.05$), and the levels of negative log-likelihood ($F(4, 36) = 110.24$, $p < 0.05$). In addition, adjusting for multiple comparisons using the Holm-Bonferroni method indicated significant differences between the normal condition ($K = 0$) and the other unusual variance conditions ($K = -8, -4, 4, 8$) in the frequencies of abnormal movements (all $p < 0.05$). Significant differences in the levels of negative log-likelihood were indicated between all pairs (all $p < 0.05$). These indicate unusual sensory variance led to unusual levels of precision-weighted prediction error, which may directly affect the perceptual processes of the robot using a prediction error minimization mechanism, thereby leading to reduction of behavioral performance.

## 4    Discussion

In this study we tested the hypothesis that aberrant sensory precision (inverse variance) causes behavioral rigidity, a core autistic behavior. In particular, using a prediction error minimization mechanism, we investigated the effects of increased and decreased sensory variance on adaptive behaviors. We conducted experiments based on a ball interaction between a humanoid robot and a human experimenter, where the robot was required to recognize situations determined by the experimenter. Although the robot with the normal network flexibly recognized situation changes and generated appropriate interactive behaviors, both increased and decreased sensory variance (inverse precision) led to seemingly similar abnormal behaviors resulting from resistance to change, such as freezing and inappropriate repetitive behavior. However, the analysis aiming to discriminate between the mechanisms underlying similar abnormal behaviors induced by the unusual variance conditions shows there were significant differences between the network-level processes underlying the symptoms and the levels of precision-weighted prediction error signals the robot experienced. Specifically, increased sensory variance resulted in disregarding any error signals, leading to invariability of intentional state, while decreased sensory variance caused an excessive response to error signals, leading to incorrect intention change and its subsequent fixation.

Our results demonstrate that increased sensory precision (decreased sensory variance) can lead to the behavioral rigidity characteristic of ASD, supporting the system-level accounts that consider increased sensory precision as the core cognitive trait of individuals with ASD[21–23, 25]. Within a theoretical study, abnormal behavioral patterns and resistance to change in individuals with ASD were proposed as strategies to provide a reassuring sense of predictive success in a world otherwise filled with error[21]. This indicates that precision-weighted prediction errors should be reduced to some extent while generating inflexible behavior. However, in our experiment, error signals could be even larger when the robot generated inflexible behavior than they were before the robot's intention was fixed. The symptoms observed in the experiment might be understood as consequences of a suboptimal solution of prediction error minimization rather than consequences of successfully reducing the sense of prediction error. However, the difference might be explained by the simplicity of our experimental setting. For example, in the experiment, the visual input to the robot was only from an external cause (a ball), but if visual inputs from internal

causes, such as the movements of the arms, were also considered, the robot might generate characteristic behaviors aiming to minimize the total error signal from the two causes by actively changing the internal causes of vision inputs. This consideration of visual inputs from internal causes might also lead to different effects on the robot's behaviors under increased or decreased sensory precision conditions.

Recently, problems with flexible adjustment of reactions to sensory states in response to volatile environments have been suggested to be associated with psychiatric disorders[25, 38]. From previous studies, not only the unusual level of reactions but also the unusual context-sensitive adjustment of reactions, such as adaptation of precision weighting of prediction errors, might explain psychiatric symptoms. In particular, a recent empirical study indicated that autistic perception may be associated with over-estimated volatility of the sensory environment, with less distinction between reactions to unexpected and expected situations[52]. In this study, sensory precision was persistently increased or decreased, indicating that influences of its context-dependent adjustment on behavioral flexibility were not considered. Future study into the effects of unusual adaptation of the precision weighting of prediction errors may facilitate understanding of finer mechanisms underlying unusual reactions to volatile environments by people with ASD. In addition, investigations will be needed of the effects of aberrant sensory precision on learning and how aberrant sensory precision can be generated through development and learning.

Our study extends attempts to understand cognitive processes underlying autistic behavior by using computational models. As a part of these attempts, Rosenberg et al.[37] conducted neural network simulations confirming that peculiarities of vision in ASD can be induced by an altered divisive normalization. Another study by Barakova et al.[33] associated poor motor skills in ASD with the poor goal-directed movements of a physical mobile robot induced by a deficit in temporal visuo-proprioceptive sensory integration. We have confirmed that aberrant sensory precision can induce behavioral rigidity, utilizing a humanoid robot controlled by a recurrent neural network model. The behavioral abnormality

was observed through a real-time human-robot interaction, where the robot was required to flexibly recognize changing environments. In such uncertain and unpredictable situations, reduced cognitive flexibility of individuals with ASD has been generally reported[5, 6]. Furthermore, this study demonstrated the generation of dysfunction in intentional control (i.e., executive dysfunction) caused by aberrant sensory precision, clarifying the direct relationship between distinct proposed cognitive abnormalities in ASD[12, 21].

Our results provide the perspective that we could consider autistic behavior as being the result of a phenomenon generally observed in natural systems. Specifically, the process leading to the qualitative shift of network behavior in the decreased sensory variance (increased sensory precision) condition may be similar to critical transitions, which are abrupt behavioral shifts observed in natural dynamical systems, including the climate, ecosystem, and cells' signaling pathways[53–55]. Critical transitions are suggested to have characteristic early warning signals, such as the slowing down of changes in a system (critical slowing down) and back-and-forth switches between states in response to relatively large impacts (flickering), although they can also occur suddenly due to a large external impact on the system[55, 56]. These characteristic phenomena were observed in network behavior in the decreased sensory variance condition. This suggests that some types of behavioral rigidity and resistance to change might result from a critical transition in the hierarchical predictive control system attributed to excessive sensory prediction errors. This perspective might be implicative because pathophysiological experiments have demonstrated that dynamical features of network behavior in epileptic seizures, which relatively high numbers of individuals with ASD experience[57], are very similar to the process of critical transition[58, 59].

Finally, findings from this study also provide an implication for clinical studies aiming to classify the different types of inflexible behavior observed in ASD or to understand differences between the behavioral abnormalities observed in ASD and other psychiatric disorders, such as obsessive-compulsive disorder and schizophrenia. Our results show that seem-

ingly similar inflexible behaviors can result from different network-level processes, and also abnormalities of network-level processes may not necessarily lead to external alterations of behavior. This indicates that measurements and classifications of behavioral abnormalities based on external observation might be confusing and create difficulties in terms of understanding their etiology as broadly described in psychiatry[14]. However, our findings also indicate that symptoms induced by increased or decreased sensory precision were substantially different in terms of the levels of prediction error signals while generating abnormal behaviors, suggesting there might be differences in the internal experiences of individuals. Therefore, measurements and classifications of both the internal experiences of patients and neural activities coding prediction error signals in the biological brain could be useful to facilitate understandings of heterogeneous behavioral rigidity in psychiatric disorders[14]. Future studies may be able to track parameters associated with those properties underlying disrupted adaptive behavior in animal models and humans, and should compare the robot model with clinical case studies.

## 5 Author summary

To behave flexibly in the uncertain world, the brain should predict whether incoming sensory information is signal or noise. Unusual estimation of the reliability (or precision) of sensory information is expected to cause unusual reactions to sensory stimuli and disabilities, especially in the complex social context. This is the basic idea of one theoretical accounting for autistic symptoms based on predictive coding, a brain principle explaining how people interact with the world as a process of minimizing the error between predictions about future states and actual sensory inputs. To bridge the gap between theoretical study and clinical observation, we performed a neuro-robotics simulation of behavioral alterations induced by over- or under- estimated sensory precision, using a recurrent neural network model based on a precision-weighted prediction error minimization mechanism. Through experiments of human–robot

adaptive interaction, we found that both increased and decreased sensory precision led to autistic-like behaviors of the robot, characterized by resistance to change, such as freezing and inappropriate repetitive behavior. In addition, the behaviors induced in the robot by increased and decreased sensory precision were different in the network-level processes and the levels of prediction error signals the robot experienced. This might indicate that the same overt symptoms could arise from distinct computational mechanisms. The results might contribute to a more transdiagnostic understanding of behavioral inflexibility.

## References

[1] Simmons, D. R. et al. Vision in autism spectrum disorders. Vision Res. 49, 2705–2739 (2009).

[2] Gowen, E. & Hamilton, A. Motor abilities in autism: A review using a computational context. J. Autism Dev. Disord. 43, 323–344 (2013).

[3] Baron Cohen, S. Theory of mind and autism: a review. Russell J. Bertrand Russell Arch. 23, 169–184 (2001).

[4] American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders 5th edn (American Psychiatric Association, 2013).

[5] Leekam, S. R., Prior, M. R. & Uljarevic, M. Restricted and repetitive behaviors in autism spectrum disorders: a review of research in the last decade. Psychol. Bull. 137, 562–593 (2011).

[6] Poljac, E. & Bekkering, H. A review of intentional and cognitive control in autism. Front. Psychol. 3, 1–15 (2012).

[7] Poljac, E., Hoofs, V., Princen, M. M. & Poljac, E. Understanding Behavioural Rigidity in Autism Spectrum Conditions: The Role of Intentional Control. J.Autism Dev.Disord. 0, (2017).

[8] Zandt, F., Prior, M. & Kyrios, M. Repetitive behaviour in children with high functioning autism

and obsessive compulsive disorder. J. Autism Dev. Disord. 37, 251–259 (2007).

[9] Lewis, M. & Kim, S. J. The pathophysiology of restricted repetitive behavior. J. Neurodev. Disord. 1, 114–132 (2009).

[10] Lord, C. & Jones, R. M. Re-thinking the classification of autism spectrum disorders. J. Child Psychol. Psychiatry 53, 490–509 (2012).

[11] Duncan, A., Ph, D., Kreiger, A., Buja, A. & Lord, C. Subcategories of restricted and repetitive behaviors in children with ASD. J. Autism Dev. Disord. 43, 1287–1297 (2013).

[12] Hill, E. L. Executive dysfunction in autism. Trends Cogn. Sci. 8, 26–32 (2004).

[13] Happé, F. & Frith, U. The weak coherence account: Detail-focused cognitive style in autism spectrum disorders. J. Autism Dev. Disord. 36, 5–25 (2006).

[14] Redish, A. D. & Gordon, J. A. (eds) Computational Psychiatry: New Perspectives on Mental Illness (MIT Press, 2016).

[15] Bar, M. The proactive brain: using analogies and associations to generate predictions. Trends Cogn. Sci. 11, 280–289 (2007).

[16] Den Ouden, H. E. M., Kok, P. & de Lange, F. P. How prediction errors shape perception, attention, and motivation. Front. Psychol. 3, 1–12 (2012).

[17] Friston, K. J., Daunizeau, J., Kilner, J. & Kiebel, S. J. Action and behavior: A free-energy formulation. Biol. Cybern. 102, 227–260 (2010).

[18] Pellicano, E. & Burr, D. When the world becomes 'too real': A Bayesian explanation of autistic perception. Trends Cogn. Sci. 16, 504–510 (2012).

[19] Friston, K. J., Lawson, R. & Frith, C. D. On hyperpriors and hypopriors: Comment on Pellicano and Burr. Trends Cogn. Sci. 17, 1 (2013).

[20] van Boxtel, J. J. A. & Lu, H. A predictive coding perspective on autism spectrum disorders. Front. Psychol. 4, 1–3 (2013).

[21] Van de Cruys, S. et al. Precise minds in uncertain worlds: predictive coding in autism. Psychol. Rev. 121, 649–75 (2014).

[22] Lawson, R. P., Rees, G. & Friston, K. J. An aberrant precision account of autism. Front. Hum. Neurosci. 8, 302 (2014).

[23] Van de Cruys, S., Van der Hallen, R. & Wagemans, J. Disentangling signal and noise in autism spectrum disorder. Brain Cogn. 112, 78–83 (2017).

[24] van Schalkwyk, G. I., Volkmar, F. R. & Corlett, P. R. A Predictive Coding Account of Psychotic Symptoms in Autism Spectrum Disorder. J. Autism Dev. Disord. 47, 1323–1340 (2017).

[25] Palmer, C. J., Lawson, R. P. & Hohwy, J. Bayesian approaches to autism: Towards volatility, action, and behavior. Psychol. Bull. 143, 521–542 (2017).

[26] Montague, P.R., Dolan, R.J., Friston, K.J. & Dayan, P. Computational psychiatry. Trends Cogn. Sci. 16, 72–80 (2012).

[27] Friston, K. J., Stephan, K. E., Montague, R. & Dolan, R. J. Computational psychiatry: The brain as a phantastic organ. The Lancet Psychiatry 1, 148–158 (2014).

[28] Adams, R. A., Huys, Q. J. M. & Roiser, J. P. Computational Psychiatry: towards a mathematically informed understanding of mental illness. J. Neurol. Neurosurg. Psychiatry jnnp-2015-310737- (2015).

[29] Huys, Q. J. M., Maia, T. V & Frank, M. J. Computational psychiatry as a bridge from neuroscience to clinical applications. Nat Neurosci 19, 404–413 (2016).

[30] Teufel, C. & Fletcher, P. C. The promises and pitfalls of applying computational models

to neurological and psychiatric disorders. Brain 139, 2600–2608 (2016).

[31] O'Loughlin, C. & Thagard, P. Autism and coherence: A computational model. Mind Lang. 15, 375–392 (2000).

[32] Diwadkar, V. A. et al. Impaired associative learning in schizophrenia: Behavioral and computational studies. Cogn. Neurodyn. 2, 207–219 (2008).

[33] Barakova, E. I. & Chonnaparamutt, W. Timing sensory integration: Robot simulation of autistic behavior. IEEE Robot. Autom. Mag. 16, 51–58 (2009).

[34] Yamashita, Y. & Tani, J. Spontaneous prediction error generation in schizophrenia. PLoS One 7, (2012).

[35] Brown, H., Adams, R. A., Parees, I., Edwards, M. & Friston, K. Active inference, sensory attenuation and illusions. Cogn. Process. 14, 411–427 (2013).

[36] Krichmar, J. L. A neurorobotic platform to test the influence of neuromodulatory signaling on anxious and curious behavior. Front. Neurorobot. 7, 1–17 (2013).

[37] Rosenberg, A., Patterson, J. S. & Angelaki, D. E. A computational perspective on autism. Proc. Natl. Acad. Sci. U. S. A. 112, 9158–9165 (2015).

[38] Powers, A. R., Mathys, C. & Corlett, P. R. Pavlovian conditioning—induced hallucinations result from overweighting of perceptual priors. Science (80-. ). 357, 596–600 (2017).

[39] Smith, L. & Gasser, M. The development of embodied cognition: six lessons from babies. Artif. Life 11, 13–29 (2005).

[40] Asada, M. et al. Cognitive Developmental Robotics: A Survey. IEEE Trans. Auton. Ment. Dev. 1, 12–34 (2009).

[41] Elman, J. L. Finding structure in time. Cogn. Sci. 14, 179–211 (1990).

[42] Marocco, D., Cangelosi, A., Fischer, K. & Belpaeme, T. Grounding action words in the sensorimotor interaction with the world: Experiments with a simulated icub humanoid robot. Front. Neurorobot. 4, 1–15 (2010).

[43] Alnajjar, F., Yamashita, Y. & Tani, J. The hierarchical and functional connectivity of higher-order cognitive mechanisms: Neurorobotic model to investigate the stability and flexibility of working memory. Front. Neurorobot. 7, 1–13 (2013).

[44] Murata, S., Namikawa, J., Arie, H., Sugano, S. & Tani, J. Learning to reproduce fluctuating time series by inferring their time-dependent stochastic properties: Application in Robot learning via tutoring. IEEE Trans. Auton. Ment. Dev. 5, 298–310 (2013).

[45] Tani, J., Ito, M. & Sugita, Y. Self-organization of distributedly represented multiple behavior schemata in a mirror system: Reviews of robot experiments using RNNPB. Neural Networks 17, 1273–1289 (2004).

[46] Ito, M., Noda, K., Hoshino, Y. & Tani, J. Dynamic and interactive generation of object handling behaviors by a small humanoid robot using a dynamic neural network model. Neural Networks 19, 323–337 (2006).

[47] Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by back-propagating errors. Nature 323, 533–536 (1986).

[48] Chen, Y. et al. Emergence of Interactive Behaviors between Two Robots by Prediction Error Minimization Mechanism. 2016 Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob), Cergy-Pontoise, 302–307 (2016).

[49] Eagleman, D. M. & Sejnowski, T. J. Motion integration and postdiction in visual awareness. Science 287, 2036–2038, 10.1126/science.287.5460.2036 (2000).

[50] Shimojo, S. Postdiction: Its implications on visual awareness, hindsight, and sense of agency. Front. Psychol. 5, 1–19 (2014).

[51] Holm, S. A simple sequentially rejective multiple test procedure. Scandinavian Journal of Statistics 6, 65–70 (1979).

[52] Lawson, R. P., Mathys, C. & Rees, G. Adults with autism overestimate the volatility of the sensory environment. Nat. Neurosci. 20, 1293–1299 (2017).

[53] May, R. M. Thresholds and breakpoints in ecosystems with a multiplicity of stable states. Nature 269, 471–477 (1977).

[54] Lenton, T. M. et al. Using GENIE to study a tipping point in the climate system. Philos. Trans. R. Soc. A Math. Phys. Eng. Sci. 367, 871–884 (2009).

[55] Scheffer, M. et al. Early-warning signals for critical transitions. Nature 461, 53–59 (2009).

[56] Scheffer, M. et al. Anticipating Critical Transitions. Science (80-. ). 338, 344–348 (2012).

[57] Bolton, P. F. et al. Epilepsy in autism: features and correlates. Br. J. Psychiatry 198, 289–294 (2011).

[58] Jiruska, P. et al. High-Frequency Network Activity, Global Increase in Neuronal Activity, and Synchrony Expansion Precede Epileptic Seizures In Vitro. J. Neurosci. 30, 5690–5701 (2010).

[59] Kramer, M. a. et al. Human seizures self-terminate across spatial scales via a critical transition. Proc. Natl. Acad. Sci. 109, 21116–21121 (2012).

## 6    Data availability

The data that support the findings of this study are available from the corresponding author on reasonable request.

## 7    Author Contributions

HI, SM, YC, YY, JT and TO conceived the research topic, designed the experiment, and wrote the paper. HI performed the experiment and analyzed the data.

## 8    Acknowledgments