Review article

# Active inference leads to Bayesian neurophysiology

Takuya Isomura

*Brain Intelligence Theory Unit, RIKEN Center for Brain Science, 2-1 Hirosawa, Wako, Saitama 351-0198, Japan*

ABSTRACT

The neuronal substrates that implement the free-energy principle and ensuing active inference at the neuron and synapse level have not been fully elucidated. This Review considers possible neuronal substrates underlying the principle. First, the foundations of the free-energy principle are introduced, and then its ability to empirically explain various brain functions and psychological and biological phenomena in terms of Bayesian inference is described. Mathematically, the dynamics of neural activity and plasticity that minimise a cost function can be cast as performing Bayesian inference that minimises variational free energy. This equivalence licenses the adoption of the free-energy principle as a universal characterisation of neural networks. Further, the neural network structure itself represents a generative model under which an agent operates. A virtue of this perspective is that it enables the formal association of neural network properties with prior beliefs that regulate inference and learning. The possible neuronal substrates that implement prior and posterior beliefs and how to empirically examine the theory are discussed. This perspective renders brain activity explainable, leading to a deeper understanding of the neuronal mechanisms underlying basic psychology and psychiatric disorders in terms of an implicit generative model.

## 1. Introduction

Biological organisms optimise their perceptions and actions to predict the external milieu and achieve preferred outcomes. The mathematical characterisation of such a sentient behaviour is crucial to better understand basic neuropsychology and psychiatric disorders. In neuroscience, two major theoretical modelling approaches exist. One employs dynamical systems derived from the physiological knowledge (or laws) of neural networks. This approach employs simple differential equations that are interpretable in terms of standard neural network architecture, for example, in the form of reservoir networks (Sussillo and Abbott, 2009; Laje and Buonomano, 2013). The other approach is the view of the brain as an agent that performs variational Bayesian inference, as considered in the Bayesian brain hypothesis (Knill and Pouget, 2004; Doya et al., 2007) and the free-energy principle and ensuing active inference (Friston et al., 2006; Friston, 2010). This perspective enables the deployment of the concept of statistical inference—established in machine learning and applied mathematics (Bishop, 2006)—to encapsulate brain functions and biological characteristics. Nevertheless, the correspondence between these two approaches is not completely understood.

Therefore, this Review aims to introduce recent progress in theoretical neurobiology, demonstrating that a class of standard neural networks performs variational Bayesian inference under the form of a generative model (Friston, 2013, 2019; Parr et al., 2020; Isomura and Friston, 2020; Isomura et al., 2022). These works demonstrate that standard neural networks—comprising biologically plausible neural activity and plasticity models—can perform Bayes optimal inference, learning, control, and planning in a self-organising manner. This notion enables to render any neural activity and plasticity 'explainable' in terms of Bayesian inference. It further enables the estimation of implicit prior beliefs—under which the neural network of the agent operates. This correspondence is the key to increase the testability of neuronal substrates underlying the free-energy principle. The remainder of the paper describes the foundations of the free-energy principle, formal correspondence to standard neural networks, possible neuronal substrates, and perspectives with regard to computational psychiatry.

## 2. Overview of the free-energy principle

The free-energy principle, proposed by Friston, explains brain functions in a unified mathematical way under the framework of Bayesian inference (Friston et al., 2006; Friston, 2010). Historically, the basis of this principle is the concept of unconscious inference (Helmholtz, 1925), as espoused by the 19th century physicist and physician Helmholtz. He proposed that humans make unconscious inferences to supplement

insufficient information given the incompleteness of sensory data and that this process underwrites perception. The brain is thus viewed as an agent that unconsciously infers the hidden dynamics underlying sensory inputs. In addition to the conceptual framework, implementations of such inferences have been considered in neuroscience and machine learning literature (Dayan et al., 1995). In particular, predictive coding is a popular scheme in which the brain is considered to update its internal representation by minimising prediction error; this has been applied to model information processing in the visual cortex (Rao and Ballard, 1999) and other brain areas. The free-energy principle formularises unconscious inference in terms of Bayesian inference of external milieu states. Further, the same principle provides a plausible explanation for adaptive behavioural control and decision-making, which is referred to as active inference (Friston et al., 2011, 2016; Friston et al., 2017a). Fig. 1 illustrates how perceptual and active inference operates under the free-energy principle.

The free-energy principle posits that minimisation of input surprise—i.e., the improbability of sensory inputs—is a law governing biological organisms, thereby providing a universal characterisation of their properties. Surprise is defined as the negative log probability of sensory inputs. In this regard, the surprise is large when receiving an unexpected input, and its minimisation indicates better adaptation to a given environment. However, it should be noted that this surprise is a statistically defined measure that is conceptually distinct to the conscious experience of feeling surprised. Given that the calculation of surprise is intractable for neural networks because it involves the marginalisation or integral of the joint probability distribution placed in the logarithm function, neural networks are thought to evaluate an upper bound of input surprise—referred to as variational free-energy—as a tractable proxy to compute surprise. The term 'free-energy principle' is derived from this concept. Thus, under this framework, neural activity and synaptic strengths are updated and actions are generated to minimise variational free energy. This property is reminiscent of Le

Chatelier's principle in thermodynamics and chemistry. Accordingly, a neural network self-organises to perform variational Bayesian inference of the external milieu states, underwriting various brain functions.

Technically, variational Bayesian inference is the process of updating the prior belief about the external milieu states to the corresponding posterior belief, based on a sequence of sensory inputs or observations ($o$). Such inference is based on a (hierarchical) generative model that mechanically expresses how external milieu states generate sensory inputs (Friston, 2008). Hereafter, the external milieu states ($\vartheta$) are defined as a set of hidden states ($s$), action (or decision) of agent ($\delta$), parameters ($\theta$), and hyper-parameters ($\lambda$), denoted as $\vartheta = \{s, \delta, \theta, \lambda\}$ (please note the difference between $\vartheta$ and $\theta$; here, one may employ policy $\pi$ instead of a sequence of actions $\delta$ to construct $\vartheta$). For example, when the external milieu is a discrete state space, it is expressed in the form of a partially observed Markov decision process (Friston et al., 2017a). In many cases, the posterior expectation $\boldsymbol{\vartheta}$ (i.e., an estimator of $\vartheta$ based on observations, or its counterpart) is sufficient to approximate the posterior belief. Hence, variational free energy *(F)* is given as a function of $o$ and $\boldsymbol{\vartheta}$ as follows:

$$F(o, \boldsymbol{\vartheta}) = [\text{Prediction error}] + [\text{Complexity}] \qquad (1)$$

Variational free energy comprises the sum of prediction error and complexity. Prediction error measures the degree to which predictions of states and inputs differ from their actual values, which is reduced to a widely used mean squared error when the background noise is considered as Gaussian (Friston, 2008). Complexity refers to the difference between the prior and posterior distributions, which is usually evaluated using Kullback-Leibler divergence. This term plays the role of regularising the posterior distribution to prevent it from being too far from the corresponding prior distribution. Minimisation of *F* with respect to components of $\boldsymbol{\vartheta}$ via a gradient descent rule optimises the posterior belief:
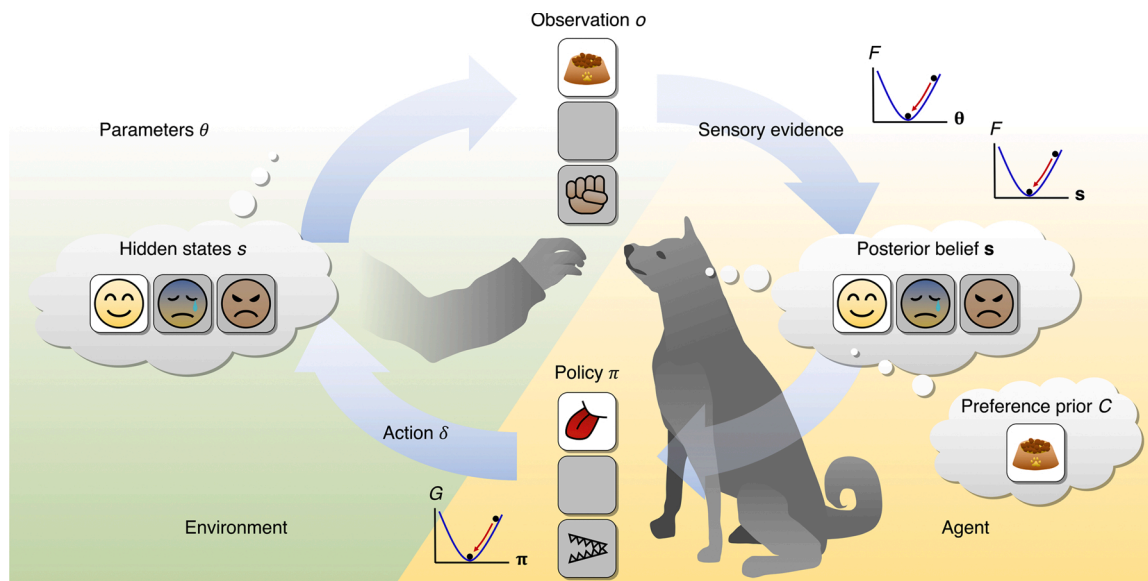


**Fig. 1. Modelled situation of perceptual and active inference under the free-energy principle.** In this schematic, we consider a dog as an agent exposed to the external world. The dog has a master, and the outcome (*o*) changes depending on the hidden (mental) states of the master (*s*). The dog employs the posterior belief of hidden states (**s**). By inferring the relationship between the hidden states and outcomes, the dog can infer what has occurred in the external (hidden state) milieu and predict what will occur next. This is attained by the Bayesian belief updating, by which the dog finds an optimised hidden state representation that minimises variational free energy *F*. The generative model that determines how the hidden states generate the observations is parameterised by *θ*. Thus, for accurate inference, the agent dog further needs to minimise the same variational free energy *F* with respect to the parameters *θ*. By learning the posterior belief of *θ*, the agent can attain a plausible generative model. The same principle can be applied to behavioural control. Under the active inference, the prior preference regarding future outcomes *C* determines the action or decision of the agent. The agent selects an action to best match future outcomes with the prior preference. Essentially, the agent needs to compute expected free energy *G* accumulated in the future. Here, *π* denotes a vector of policies—parameters determining a sequence of actions—that represents the rule for selecting behaviour depending on the state posterior. For instance, if *π* is 1, the dog selects licking behaviour. By selecting a policy that minimise the expected free energy, the dog attains a treat.

$$\dot{\boldsymbol{\vartheta}}_i \propto -\frac{\partial F}{\partial \boldsymbol{\vartheta}_i} \qquad (2)$$

where $\boldsymbol{\vartheta}_i$ indicates the *i*th component of $\boldsymbol{\vartheta}$. This update of $\boldsymbol{\vartheta}$ attains a fixed point (i.e., $\boldsymbol{\vartheta}$ that gives $\dot{\boldsymbol{\vartheta}} = 0$), which means that the representation and behaviour are Bayes optimal. Thus, the free-energy principle provides unified explanations of inference (i.e., optimisation of *s*), learning (that of **θ**), adaptive behavioural control (that of **δ**), predictions of *o* and *s* in the future, and the ensuing planning to minimise the risk associated with future outcomes.

Under this framework, the generative model—which is also referred to as the internal model (Dayan et al., 1995; George and Hawkins, 2009)—represents the agent's hypothesis about how external milieu states generate sensory inputs, whereby perceptual learning is cast as the optimisation of the generative model. The generative model self-organises to match the actual generative process of sensory inputs; thereby, the activity of the agent's neural network can accurately infer external milieu states and predict subsequent states and inputs. This yields an information representation referred to as predictive coding (Rao and Ballard, 1999; Friston, 2005). Related to this, the optimisation under some sparse priors yields efficient internal representation, referred to as sparse coding (Olshausen and Filed, 1996).

A virtue of the free-energy principle is that it applies Bayesian inference to explain the optimisation of action and planning, i.e., active inference (Friston et al., 2011, 2016; Friston et al., 2017a). When an agent returns feedback responses to the external milieu, the generative process—and the ensuing surprise—become a function of the agent's action. Thus, the agent generates an action to minimise the expected value of variational free energy in the future, which is referred to as expected free energy (*G*) (Friston et al., 2017a), thereby rendering the actual observation closer to its preferred (i.e., predicted) outcome. This preference is characterised in the form of the preference prior (*C*). For example, the agent (dog) selects a behaviour that minimises the expected free energy to obtain food (Fig. 1).

Active inference occurs when the agent receives sensory inputs differ from what it predicted. For instance, this happens when the agent employs a generative model that differs from the generative process of the external milieu. An action is thus generated to render the external milieu generative process closer to the generative model that the agent employs (Friston et al., 2011). An example of this scenario is birdsong (Kiebel et al., 2008; Friston and Frith, 2015a, b). Once the agent (in this case, a bird) has learnt the state of hearing the song of other birds, the presence of the song minimises surprise. Thus, when the agent does not hear any song, it tries to hear a song, such as by singing itself or by seeking conspecifics, because the absence of the song yields a large surprise. As a consequence of action generation, the agent receives its own prediction (that is, predicted song) as the actual sensory inputs, thus minimising surprise. An extension of this model to a mixture of experts (Wolpert and Kawato, 1998) allows an agent employing multiple generative models to learn to predict and imitate songs of several different birds (Isomura et al., 2019). Note that the bird may re-adapt to the absence of songs prior to action generation. In short, surprise minimisation occurs in two ways: the agent's internal states may approach the external milieu states, or the agent's action may make the external milieu states get closer to the internal states. The balance between the learning rate and the threshold for action generation determines whether learning or action generation may occur.

Active inference also underwrites planning (Friston et al., 2016, 2017a). Planning corresponds to the selection of behavioural policies to minimise future uncertainty; i.e., planning as inference (Attias, 2003; Botvinick and Toussaint, 2012; Maisto et al., 2015; Kaplan and Friston, 2018; Millidge, 2020). While action (*δ*) directly affects the external milieu, policy (*π*) represents future plans, or equivalently a sequence of actions, which corresponds to the parameters that determine actions. The posterior beliefs of policies are proportional to the exponential of negative expected free energies multiplied by precision. Thus, the agent computes the expected free energy associated with each policy and selects the option that provides the minimum expected free energy. Here, the prior preference about future outcomes (which involves information on rewards and punishments) characterises the shape of the expected free energy. In relation to neurophysiology, the optimisations of policies may be associated with *post-hoc* modulation of synaptic plasticity mediated by various neuromodulators, which will be discussed in the subsequent section.

In active inference, the balance of exploitation and exploration is determined by expected free energy. If a policy gives a much smaller expected free energy than others, the probability of the policy being selected is close to 1, leading to an exploitative strategy. Conversely, if all the policies give similar expected free energies, the agent randomly selects a policy, yielding an explorative behaviour. Furthermore, the precision that controls the magnitude of expected free energy is optimised by minimising variational free energy, where a higher precision renders the agent's behaviour more exploitative.

In summary, the free-energy principle is a rule that biological organisms follow. The internal representations of organisms are optimised by minimising variational free energy, rendering the brain activity to exhibit predictive coding. The same principle derives active inference, whereby behaviours are optimised to minimise the expected free energy in order to obtain preferred outcomes, providing a plausible explanation for planning and goal-directed behaviour. Thus, both perceptual and active inference contribute to the minimisation of variational free energy, accumulated over either the past or future. This framework provides a universal characterisation of the sentient behaviour of biological organisms.

## 3. Mathematical equivalence between neural networks and variational Bayes

The free-energy principle is a theory with high abstraction, and its neuronal mechanisms remain to be fully elucidated (Fig. 2), although the state and parameter posteriors are usually associated with neural activity and synaptic strengths (Friston et al., 2006; Friston, 2010) and evidence has been accumulated (Bastos et al., 2012). Generally, the distinction between a biological organism and its surrounding
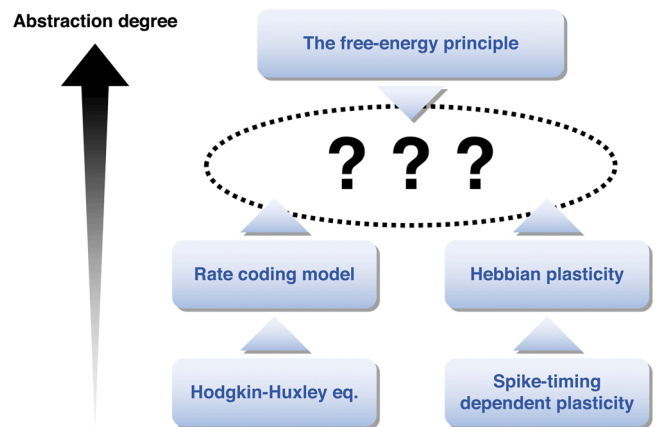


**Fig. 2. Abstraction degree of theories.** Although the universality of the Bayesian inference framework permits the empirical explanation of brain activity and behaviour in a unified manner, this property also makes it difficult to identify the neuronal substrates underlying the free-energy principle. This is because there exist many possible implementations of Bayesian inference that are equally good from the theoretical perspective. Thus, identification of the computational architecture by which neural networks implement the free-energy principle remains challenging. Conversely, experimental validations of computational models of neural activity and synaptic plasticity have been relatively established. Hence, filling the gap between the top-down principle and bottom-up physiological phenomena is crucial for identification.

environment implies the existence of a Markov blanket that statistically separates internal and external states. When a system reaches a (possibly non-equilibrium) steady state, the conditional expectation of the internal states of a biological organism is cast as parameterising posterior beliefs about the external milieu states (Friston, 2013, 2019; Parr et al., 2020). This implies that any (non-equilibrium) steady state realises some Bayesian inference. To rephrase, according to the complete class theorem, there exists at least a set of priors and Bayesian cost functions that can explain any observed behaviour of a biological agent in terms of Bayesian inference (Wald, 1947; Brown, 1981; Berger, 2013). This means that "one might not be able to experimentally refute the hypothesis that the brain acts as a Bayesian observer" (Daunizeau et al., 2010). One might think this property is problematic when designing experimental validations of the principle.

In contrast, theories of brain elements such as neurons and synapses are more well established. Neurons generate spiking activities, and the nonlinear dynamics of neuronal membrane potentials and ion channels are explained by the Hodgkin-Huxley equation (Hodgkin and Huxley, 1952). These complex dynamics can be reduced to a rate-coding model (Adrian and Zotterman, 1926), which determines neural activity as a function of sensory inputs and past network activity. The activity (i.e., firing intensity) is updated by the sum of the leak factor, synaptic inputs weighted by synaptic strengths, and firing threshold, wherein the threshold determines the mean firing level. The activity is often characterised by the sigmoid function, also known as neurometric function (Newsome et al., 1989).

Moreover, synaptic strengths exhibit plasticity depending on neural activity. Activity-dependent plasticity is governed by the timing between pre- and post-synaptic activity, referred to as spike-timing dependent plasticity (STDP) (Markram et al., 1997; Bi and Poo, 1998). STDP is observed in various brain regions both *in vivo* and *in vitro*. When neurons adopt rate coding—which is plausible considering large noise in the brain and a high robustness of rate coding against the noise (London et al., 2010)—STPD is reducible to the Hebbian plasticity rule (Clopath et al., 2010). Hebb's law (Hebb, 1949) states that synaptic strengths ($W$) are strengthened when pre- and post-synaptic neurons fire together. This is expressed as the product of pre- and post-synaptic activity: $\dot{W} \propto pre \times post$. Here, $\dot{W}$ denotes the change in synaptic strengths. This law enhances the association between (pre-synaptic) causes and (post-synaptic) consequences. Electrophysiological experiments showed that Hebbian plasticity occurs depending on the activity level (Bliss and Lømo, 1973; Malenka and Bear, 2004), spike timings (Markram et al., 1997; Bi and Poo, 1998), or burst timings (Butts et al., 2007) of pre- and post-synaptic neurons. These rules of activity and plasticity provide concise and plausible descriptions of the neural network dynamics.

These basic properties presumably remain unaltered even for neurons integrated into a large network, i.e., in the brain. Further, owing to the difference in the time scales of neural activity and plasticity, there exists a cost function that derives both neural activity and synaptic plasticity (Isomura et al., 2022). Thus, let us consider a simple neural network model wherein the internal states of a neural network ($\varphi$) comprise a set of neural activity ($x, y$), synaptic weights ($W$), and other free parameters ($\phi$) such as firing threshold factors, denoted as $\varphi = \{x, y, W, \phi\}$. The dynamics of these internal states are expressed in terms of the gradient descent on a biologically plausible cost function for the neural network $L$:

$$\dot{\varphi}_i \propto -\frac{\partial L}{\partial \varphi_i} \tag{3}$$

where $\varphi_i$ indicates the $i$th component of $\varphi$. This update rule is a canonical expression of the dynamics of standard neural networks. Here, if $\varphi_i$ is $x$, this becomes the activity rule; if $\varphi_i$ is $W$, this becomes the plasticity rule. The form of the cost function $L$ can be identified by simply taking the integral of the neural activity equation, such as the rate coding model, referred to as the reverse-engineering approach (Isomura

and Friston, 2020). Crucially, the synaptic update rule derived as the gradient descent on the same cost function $L$ has a biologically plausible form comprising Hebbian plasticity accompanied by a homeostatic plasticity term. Thus, the minimisation of cost function $L$ is sufficient to characterise the computational architecture of the standard neural network. Given this, one may refer to $L$ as complete neural network potential.

However, it remains unclear how these dynamics underwrite brain functions; although the free-energy principle provides an empirical explanation, how standard neural networks can implement the principle remains to be fully understood. In this regard, the elucidation of the formal link between the top-down principle and bottom-up physiological phenomena at the neuron and synapse level is essential (Fig. 2). Designing a network that exhibits various functions—while comprising standard activity and plasticity models—enables to better understand circuit mechanisms underlying brain functions, i.e., constructive approach.

To address this issue, recent work has advocated that any neural network cost function $L$ can be cast as variational free energy $F$ (Isomura and Friston, 2020; Isomura et al., 2022). Therefore, the gradient flow on $L$, depicted in Eq. (3), is mathematically equivalent to that on $F$, in Eq. (2). This indicates that the dynamics of the standard neural network implicitly performs variational Bayesian inference and learning. Therefore, the neural network interacting with the environment is cast as performing active inference (Fig. 3). This property owes to the complete class theorem that states that any admissible decision rule can be cast as the Bayes optimal solution under at least a pair of Bayesian cost function and prior beliefs (Wald, 1947; Brown, 1981; Berger, 2013). The form of variational free energy is characterised by the generative model under which the agent operates. Thus, this theory argues that for any given neural network that minimises $L$, there exists a generative model that satisfies

$$F(o, \boldsymbol{\vartheta}) \equiv L(o, \varphi) \tag{4}$$

Here, the neural network's internal states $\varphi$ encode or parameterise the posterior expectation of the external milieu states $\boldsymbol{\vartheta}$. Please refer to (Isomura and Friston, 2020; Isomura et al., 2022) for technical details and examples of analytically tractable neural network architectures that
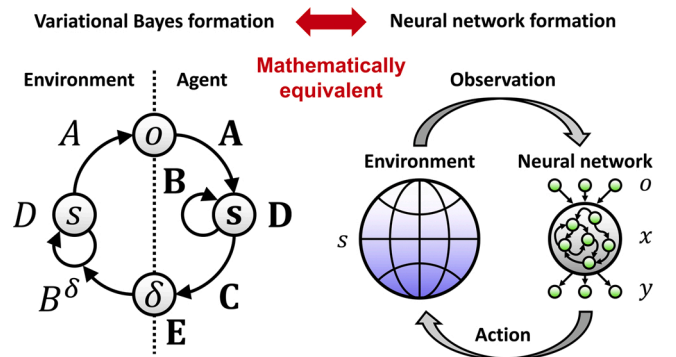


**Fig. 3. Schematic of a neural network interacting with the environment and corresponding Bayesian formation.** This schematic is adapted from Isomura et al. (2022). The environment is characterised by the dynamics of hidden states ($s$), and the hidden dynamics generate sensory inputs ($o$) to the agent. Left: variational Bayes formation. When the environment is characterised by discrete state space, it is expressed as a partially observed Markov decision process model. Here, matrices *A, B, C, D, E* indicate parameters that characterise the system, and *s, o, δ* are discrete variables. Right: neural network formation. A neural network consists of three layers: an input layer, a middle layer that comprises a recurrent neural network ($x$), and an output layer ($y$) that provides feedback responses to the environment. When neural activity and plasticity minimise the same cost function, such a standard neural network performs active inference.

satisfy Eq. (4).

Under this framework, one can assign the functional meaning of network dynamics in terms of Bayesian inference because all factors in the neural network formally correspond to quantities in Bayesian inference (Fig. 4). This means that although different notations are employed to explain neural networks and Bayesian models, they are homologous. Thus, standard neural networks implicitly employ a specific but generic form of the generative model. This further indicates that the dynamics of neural activity and synaptic plasticity occur in a manner that recapitulates the external dynamics within the neural network. This is an impactful explanation of plasticity and adaptation from the concept of free energy minimisation.

A virtue of this notion is that once the neural network quantities are interpreted to Bayesian inference, the functional meanings of internal states such as synaptic connections and firing thresholds are systematically determined in terms of quantities in Bayesian inference. For example, recurrent connections formally correspond to the state transition matrix, and feedforward connections from the middle to output layer formally correspond to the policy. These formal correspondences increase the explainability of the network architecture and dynamics. This perspective complements conventional approaches using naive neural network models, whereby it is difficult to interpret the functional meaning of synaptic connections and firing thresholds, which further leads to difficulties in optimising these network parameters.

The equivalence between neural networks and variational Bayes offers the identification of the optimal internal states to best represent the external milieu and generate actions. A gradient flow of the internal states that minimises a common cost function inevitably induces Bayes optimal representation and decision-making after sufficient training. Thus, standard neural networks can, in principle, achieve the Bayes optimal model. This framework provides an interesting perspective that covers multiple time scales in biology, from biochemical reaction to natural selection. For instance, the optimisation of a generative model can be associated with the functional and structural plasticity of synapses, development of an individual, and evolution of a species, depending on its time scale and hierarchy. In contrast, a naive neural

network with a suboptimal prior is unable to perform given tasks well. This occurs when the training is insufficient relative to the update of the prior belief. To rephrase, only neural networks with the optimal implicit prior belief attain the Bayes optimal model.

In this work, a rate-coding neuron was adopted owing to its analytical tractability. The rate-coding model can be derived as a reduction of realistic neuron models through some approximations, indicating its biological plausibility (Isomura et al., 2022). Further, the proposed framework can, in principle, be applied to spiking neuron models, although the identification of the implicit generative model for a given spiking model is a more delicate problem.

Related theories of the brain have been developed in which the brain is considered to maximise information (Linsker, 1988) or reduce redundancy (Barlow, 1961). The relation between the information maximisation (infomax) and free-energy principles has been discussed from multiple perspectives. Under the same generative model, these principles may provide different solutions owing to the difference in their cost functions (Isomura, 2018). Nevertheless, the cost function for the infomax principle can be cast as variational free energy under some generative model with prior beliefs according to the complete class theorem, although identifying the implicit generative model may not be straightforward.

In summary, a class of biologically plausible cost functions for neural networks can be cast as variational free energy. A benefit of this perspective is that it shows that any standard neural network becomes Bayes optimal through adaptation; thus, the dynamics that minimise the cost function render internal states—including neural activity, synaptic strengths, and firing threshold—as the Bayes optimal encoder and controller after sufficient training.

## 4. Possible neuronal substrates underlying active inference

The above consideration thus implies that neuronal substrates for the free-energy principle and active inference exist ubiquitously because any standard neural networks implicitly minimise variational free energy. This notion is a powerful guide to consider testable predictions.

Detailed correspondences between the canonical microcircuit in the cortex (Haeusler and Maass, 2007) and the hierarchical predictive coding model (Friston, 2008) have been investigated, and testable predictions have been provided from physiological and anatomical viewpoints (Bastos et al., 2012). For instance, the frequency difference in the neurons of the superficial and deep cortical layers suggests that the former and latter encode prediction errors and expectations, respectively. It has also been proposed that the state estimation is conducted in the hippocampus, and subsequently the expected free energy is computed in the ventral prefrontal cortex, which is used to select policies in the striatum (Friston et al., 2017a). A hierarchical computational architecture has been proposed in which the early sensory cortex employs a predictive coding model in a continuous state space and the higher areas adopt a partially observed Markov decision process model in a discrete state space as generative model (Friston et al., 2017b). The recent development of measurement technology will enable the investigation of the circuit architecture and verification of these predictions.

Experimental work has shown that neuronal ensembles in layers 2/3 and 5 of the parietal lobe of rodents encode posterior beliefs about hidden sound cues, which are used to reach a goal in an uncertain environment (Funamizu et al., 2016). Some zebrafish employ neuronal populations encoding state prediction errors, which facilitate fish to reach a goal (Torigoe et al., 2021). These observations suggest that these neurons perform Bayesian (active) inference associated with goal-directed behaviour. Further, according to the free-energy principle, neural circuits minimise variational free energy. Reduction of free energy was observed using *in vitro* neural networks at the neural circuit level, wherein *in vitro* networks assimilating sensory inputs self-organised to encode the state posteriors and reduced free energy (Isomura et al., 2015; Isomura and Friston, 2018).
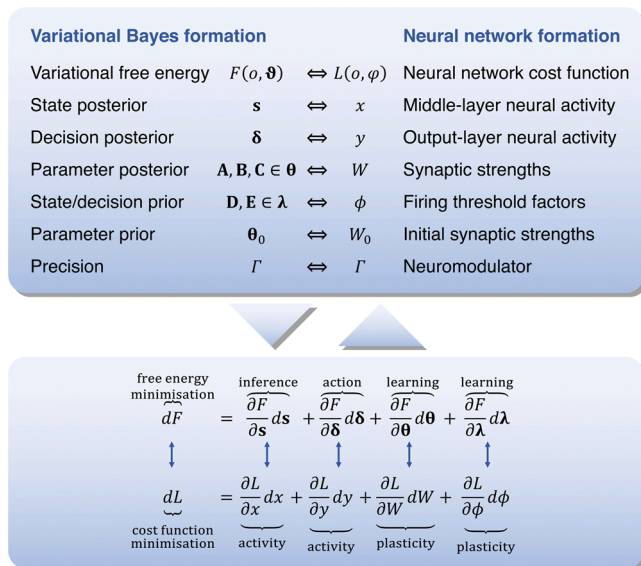


| Variational Bayes formation | | | Neural network formation |
|---|---|---|---|
| Variational free energy | $F(o, \vartheta)$ | $\Leftrightarrow$ $L(o, \varphi)$ | Neural network cost function |
| State posterior | $\mathbf{s}$ | $\Leftrightarrow$ $x$ | Middle-layer neural activity |
| Decision posterior | $\boldsymbol{\delta}$ | $\Leftrightarrow$ $y$ | Output-layer neural activity |
| Parameter posterior | $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \boldsymbol{\theta}$ | $\Leftrightarrow$ $W$ | Synaptic strengths |
| State/decision prior | $\mathbf{D}, \mathbf{E} \in \boldsymbol{\lambda}$ | $\Leftrightarrow$ $\phi$ | Firing threshold factors |
| Parameter prior | $\boldsymbol{\theta}_0$ | $\Leftrightarrow$ $W_0$ | Initial synaptic strengths |
| Precision | $\Gamma$ | $\Leftrightarrow$ $\Gamma$ | Neuromodulator |

$$\underset{\substack{\text{free energy}\\\text{minimisation}}}{dF} = \overset{\text{inference}}{\frac{\partial F}{\partial \mathbf{s}} d\mathbf{s}} + \overset{\text{action}}{\frac{\partial F}{\partial \boldsymbol{\delta}} d\boldsymbol{\delta}} + \overset{\text{learning}}{\frac{\partial F}{\partial \boldsymbol{\theta}} d\boldsymbol{\theta}} + \overset{\text{learning}}{\frac{\partial F}{\partial \boldsymbol{\lambda}} d\boldsymbol{\lambda}}$$

$$\underset{\substack{\text{cost function}\\\text{minimisation}}}{dL} = \underset{\text{activity}}{\frac{\partial L}{\partial x} dx} + \underset{\text{activity}}{\frac{\partial L}{\partial y} dy} + \underset{\text{plasticity}}{\frac{\partial L}{\partial W} dW} + \underset{\text{plasticity}}{\frac{\partial L}{\partial \phi} d\phi}$$

**Fig. 4. Correspondence of variational Bayes formation and neural network formation.** Components of $L$ and $F$ exhibit a formal correspondence. Neural activity ($x$, $y$) corresponds to the state or decision posterior ($\mathbf{s}$, $\boldsymbol{\delta}$), synaptic strengths ($W$) correspond to the parameter posterior ($\mathbf{A}$, $\mathbf{B}$, $\mathbf{C}$), firing threshold factor ($\phi$) corresponds to the state or decision prior ($\mathbf{D}$, $\mathbf{E}$), etc. Precision ($\Gamma$) may encode the information of risk, attention, gain, etc. They are summarised in a correspondence between components of the differential form of $L$ and $F$, as depicted in the lower panel.

Neuromodulators are possible neuronal substrates that implement implicit prior beliefs that regulate inference and learning (Doya, 2002; Parr and Friston, 2017). Recent studies have shown that various modulatory factors change the magnitude and sign of Hebbian plasticity in a different manner, which yield various associative functions (Pawlak et al., 2010; Frémaux and Gerstner, 2016; Kuśmierz et al., 2017). These are explained as the product of three factors (i.e., pre- and post-synaptic neural activity and modulator), which is referred to as three-factor learning rules (Fig. 5). Modulator may encode a belief—such as about precision of likelihood, prior, or posterior, learning rate, attentional filter, or risk associated with future outcomes—that nuances the form of generative model (Friston, 2008; Friston et al., 2017a; Parr and Friston, 2017; Isomura et al., 2019; Isomura et al., 2022). Various neuromodulators are known to modulate synaptic plasticity, such as dopamine (Reynolds et al., 2001; Zhang et al., 2009; Yagishita et al., 2014), noradrenaline (Salgado et al., 2012; Johansen et al., 2014), muscarine (Seol et al., 2007), and GABA (Paille et al., 2013; Hayama et al., 2013), as well as glial factors (Ben Achour and Pascual, 2010).

Prediction and planning are particularly important for active inference. However, which neural circuits realise planning and how generative models established in the sensory cortex are used for planning remain unclear. To select current actions to minimise future risk, association between different times is essential. In this regard, an important form of implicit prior belief is the expected free energy, or its reduced form, risk function. The (expected) risk is a function of the past decisions of the agent, which involves the role of reward and punishment. After observing outcomes and ensuing risk, the agent can change the evaluation of past decisions in a *post-hoc* manner, which is sufficient to achieve Bayes optimal decision-making that minimise future risk (Isomura et al., 2022). A possible substrate to enable this association is delayed modification of plasticity. Associative or Hebbian plasticity usually occurs between pre- and post-synaptic neurons based on experiences, wherein 'pre' represents the environmental information and 'post' the decisions. The post-hoc modulation enables the alteration of the association after observing its outcomes, enabling to associate past and future events through a simple mechanism. Such a post-hoc modulation of the previous associative plasticity is observed with dopamine (Yagishita et al., 2014; Wieland et al., 2015; Brzosko et al., 2017) as well as noradrenaline and serotonin (He et al., 2015).

The equivalence between neuronal circuits and variational Bayesian inference enables to the reverse-engineering of generative models and prior beliefs—that neuronal circuits implicitly employ—from empirical data (Fig. 6). This will enable the systematic estimation of the aim and internal states of animals in terms of variational Bayesian inference.
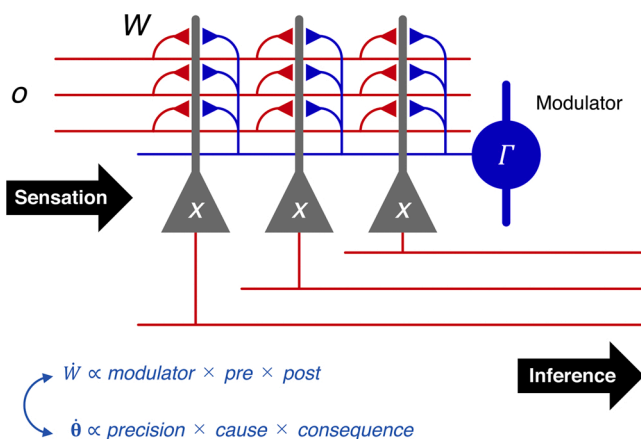
**Fig. 5. A possible minimal circuit structure realising active inference.** A group of cortical neurons exhibiting Hebbian plasticity self-organises to encode the external milieu states through associations between causes and consequences. Various modulators may regulate the plasticity, which corresponds to modulation of precision. This can realise various associative functions.
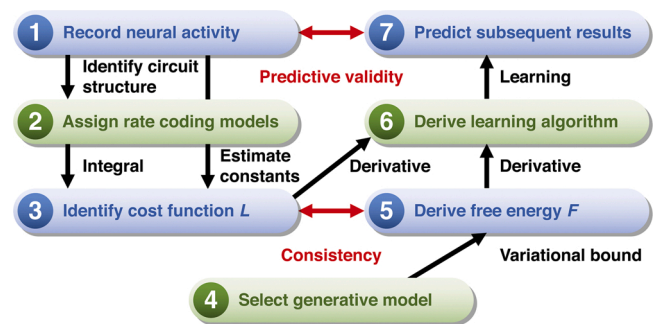
**Fig. 6. A possible procedure for experimental validation of the theory.** A network of rate-coding models is assigned to a network structure estimated based on empirical observations. Then, the biologically plausible cost function for the neural network $L$ is reverse engineered by computing the integral of the rate-coding models. The theory ensures the existence of the corresponding generative model and ensuing variational free energy $F$. From this comparison, a plausible generative model can be identified. Moreover, the derivative of $F$ (or equivalently, that of $L$) derives the plasticity rule. By computing the time evolution, one can predict the subsequent learning. Thus, in principle, the theory can predict the subsequent neural activity, plasticity, and behaviours without observation. In this manner, one can test the prediction ability of the theory.

Unlike conventional connectomics approaches that provide functional connectivity structure, the reverse-engineering approach provides a computational architecture including neural activity, plasticity, and action generation, and thus can predict the subsequent learning. This speaks to the predictive validity of the free-energy principle. The aforementioned equivalence offers predictions about what kind of plasticity and adaptation will occur in response to given experiences, wherein synaptic plasticity is cast as Bayesian belief updating to optimise the generative model. This provides a systematic procedure to investigate neuronal substrates that realise the free-energy principle using empirical data. In short, a generative model means a hypothesis that an animal employs. Once the generative model is identified, brain activity and behaviour are, in principle, predictable based on the generative model. This will be useful for designing testable predictions and examining the predictive validity of the free-energy principle.

## 5. Towards a theory of psychiatric disorders

Under the free-energy principle, the neuronal mechanisms of various psychiatric disorders are uniformly explained in terms of abnormality in inference. These disorders can be attributed to false prior beliefs that modify the perception abnormally. The posterior belief (i.e., neural activity) is determined by the sum of bottom-up sensory evidence and top-down prior beliefs weighted by precisions. Thus, the balance of integrating sensory evidence and prior beliefs is crucial for correctly perceiving the external milieu, whereas a balance break in this process leads to disordered inference and learning. From this perspective, positive symptoms in schizophrenia, such as hallucinations and delusions, can be viewed as a condition where prior beliefs are too strong (Fletcher and Frith, 2009). Autistic symptoms such as hypersensitivity to sensation and learning disabilities can be associated with overfitting or overlearning to a given environment due to abnormal precisions about prior beliefs and sensory evidence (Pellicano and Burr, 2012; Friston et al., 2014).

For instance, this property is illustrated by a force-matching illusion (Adams et al., 2013). In this task, the participant's hand is touched by a device, and thus they sense an external force. The participants are asked to reproduce the magnitude of the external force by pushing a hand with a finger of the opposite hand. Healthy people underestimate the magnitude of the detected force (i.e., self-pushing power), which results in pushing a hand with a significantly larger power than the power of the

external force, leading to the force-matching illusion. In contrast, patients with schizophrenia can perform this task more accurately. According to the free-energy principle, an attenuation of the sensory precision induces this illusion, which is usually needed to distinguish between one's own action and the external force when people control their hands. In healthy people, this sensory attenuation induces underestimation of the force generated by themselves. Conversely, this does not occur in patients with schizophrenia, leading to the unbiased inference of the magnitude of the force. However, this can be problematic in practical situations due to its close relation to hallucinations; for example, in auditory hallucinations, participants can confuse their own voice with someone else's.

The equivalence between neural network dynamics and variational Bayesian inference may enable us to systematically understand the enormous available corpus of physiological data on patients with psychiatric disorders from functional and neuronal viewpoints. The reverse-engineering approach is potentially useful as a guide to identify the neuronal mechanisms that cause these abnormalities in terms of generative models and prior beliefs and to predict the impact of abnormal neuromodulations on behaviour. This leads to characterisations of learning disabilities and biased (i.e., suboptimal) decisions in psychiatric disorder symptoms. An agent that accurately performs a given task is viewed as near Bayes optimal with appropriate prior beliefs—whereas one that fails the task is explained by false beliefs—-wherein circuit-level changes are potentially predictable from behavioural-level symptoms, and vice versa. Facilitating this may lead to the establishment of methods for early diagnosis and treatment, which can have a great social impact in terms of improvement of the quality of life for patients with psychiatric disorders.

## 6. Conclusion

In summary, the free-energy principle underwrites the optimisation of perception, learning, and action. The dynamics of standard neural networks that minimises a common cost function can be cast as Bayesian belief updating. This equivalence indicates that neural activity and synaptic plasticity accompanied by some neuromodulations are sufficient to perform Bayes optimal inference and decision-making. This implies that the recapitulation of external dynamics is an inherent feature of neural networks. Further, this theory offers the identification of the optimal parameters to encode the external milieu and perform decision-making for a wide range of neural networks. Conversely, arbitrarily selected network properties correspond to suboptimal prior beliefs, which lead to biased inference and learning and failure of prediction, thereby leading to attenuation of task performance. Thus, psychiatric disorders are associated with a balance break in the process of integrating prior beliefs and sensory evidence at the circuit level. This notion offers a universal characterisation of neural networks in terms of variational free energy minimisation (i.e., Bayesian neurophysiology) and provides testable predictions about the neuronal mechanisms underlying the sentient behaviour and its disorders. In essence, the equivalence renders brain activity explainable and affords a deeper understanding of basic neuropsychology and psychiatric disorders.

## Declaration of Competing Interest

The author reports no declarations of interest.

## Acknowledgements

## References

Adams, R.A., Stephan, K.E., Brown, H.R., Frith, C.D., Friston, K.J., 2013. The computational anatomy of psychosis. Front. Psychiatry 4, 47.

Adrian, E.D., Zotterman, Y., 1926. The impulses produced by sensory nerve-endings: Part II. The response of a Single End-Organ. J. Physiol. 61 (2), 151–171.

Attias, H., 2003. Planning by probabilistic inference. Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics.

Barlow, H., 1961. Possible principles underlying the transformations of sensory messages. In: Rosenblith, W. (Ed.), Sensory Communication. MIT Press, Cambridge MA, pp. 217–234.

Bastos, A.M., Usrey, W.M., Adams, R.A., Mangun, G.R., Fries, P., Friston, K.J., 2012. Canonical microcircuits for predictive coding. Neuron 76 (4), 695–711.

Ben Achour, S., Pascual, O., 2010. Glia: the many ways to modulate synaptic plasticity. Neurochem. Int. 57, 440–445.

Berger, J.O., 2013. Statistical Decision Theory and Bayesian Analysis. Springer Science & Business Media, Berlin.

Bi, G.Q., Poo, M.M., 1998. Synaptic modifications in cultured hippocampal neurons: dependence on spike timing, synaptic strength, and postsynaptic cell type. J. Neurosci. 18, 10464–10472.

Bishop, C.M., 2006. Pattern Recognition and Machine Learning. Springer, New York, NY, USA.

Bliss, T.V., Lømo, T., 1973. Long-lasting potentiation of synaptic transmission in the dentate area of the anaesthetized rabbit following stimulation of the perforant path. J. Physiol. 232, 331–356.

Botvinick, M., Toussaint, M., 2012. Planning as inference. Trends Cogn. Sci. 16, 485–488.

Brown, L.D., 1981. A complete class theorem for statistical problems with finite-sample spaces. Ann. Stat. 9, 1289–1300.

Brzosko, Z., Zannone, S., Schultz, W., Clopath, C., Paulsen, O., 2017. Sequential neuromodulation of Hebbian plasticity offers mechanism for effective reward-based navigation. eLife 6, e27756.

Butts, D.A., Kanold, P.O., Shatz, C.J., 2007. A burst-based "Hebbian" learning rule at retinogeniculate synapses links retinal waves to activity-dependent refinement. PLoS Biol. 5, e61.

Clopath, C., Büsing, L., Vasilaki, E., Gerstner, W., 2010. Connectivity reflects coding: a model of voltage-based STDP with homeostasis. Nat. Neurosci. 13 (3), 344–352.

Daunizeau, J., Den Ouden, H.E., Pessiglione, M., Kiebel, S.J., Stephan, K.E., Friston, K.J., 2010. Observing the observer (I): Meta-Bayesian models of learning and decision-making. PLOS One 5, e15554.

Dayan, P., Hinton, G.E., Neal, R.M., Zemel, R.S., 1995. The Helmholtz machine. Neural Comput. 7, 889–904.

Doya, K., 2002. Metalearning and neuromodulation. Neural Netw. 15, 495–506.

Doya, K., Ishii, S., Pouget, A., Rao, R.P. (Eds.), 2007. Bayesian Brain: Probabilistic Approaches to Neural Coding. MIT Press, Cambridge, MA, USA.

Fletcher, P.C., Frith, C.D., 2009. Perceiving is believing: a Bayesian approach to explaining the positive symptoms of schizophrenia. Nat. Rev. Neurosci. 10 (1), 48–58.

Frémaux, N., Gerstner, W., 2016. Neuromodulated spike-timing-dependent plasticity, and theory of three-factor learning rules. Front. Neural Circuits 9, 85.

Friston, K.J., 2005. A theory of cortical responses. Philos. Trans. R. Soc. Lond. B Biol. Sci. 360, 815–836.

Friston, K.J., 2008. Hierarchical models in the brain. PLoS Comput. Biol. 4 (11), e1000211.

Friston, K.J., 2010. The free-energy principle: a unified brain theory? Nat. Rev. Neurosci. 11, 127–138.

Friston, K.J., 2013. Life as we know it. J. R. Soc. Interface 10, 20130475.

Friston, K.J., 2019. A free energy principle for a particular physics. arXiv preprint: 1906.10184.

Friston, K.J., Frith, C.D., 2015a. Active inference, communication and hermeneutics. Cortex 68, 129–143.

Friston, K.J., Frith, C.D., 2015b. A duet for one. Conscious. Cogn. 36, 390–405.

Friston, K.J., Kilner, J., Harrison, L., 2006. A free energy principle for the brain. J. Physiol. Paris 100, 70–87.

Friston, K.J., Mattout, J., Kilner, J., 2011. Action understanding and active inference. Biol. Cybern. 104, 137–160.

Friston, K.J., Stephan, K.E., Montague, R., Dolan, R.J., 2014. Computational psychiatry: the brain as a phantastic organ. Lancet Psychiatry 1 (2), 148–158.

Friston, K.J., FitzGerald, T., Rigoli, F., Schwartenbeck, P., O'Doherty, J., Pezzulo, G., 2016. Active inference and learning. Neurosci. Biobehav. Rev. 68, 862–879.

Friston, K.J., FitzGerald, T., Rigoli, F., Schwartenbeck, P., Pezzulo, G., 2017a. Active inference: a process theory. Neural Comput. 29, 1–49.

Friston, K.J., Parr, T., de Vries, B.D., 2017b. The graphical brain: belief propagation and active inference. Netw. Neurosci. 1, 381–414.

Funamizu, A., Kuhn, B., Doya, K., 2016. Neural substrate of dynamic Bayesian inference in the cerebral cortex. Nat. Neurosci. 19 (12), 1682–1689.

George, D., Hawkins, J., 2009. Towards a mathematical theory of cortical micro-circuits. PLoS Comput. Biol. 5 (10), e1000532.

Haeusler, S., Maass, W., 2007. A statistical analysis of information-processing properties of lamina-specific cortical microcircuit models. Cereb. Cortex 17 (1), 149–162.

Hayama, T., Noguchi, J., Watanabe, S., Takahashi, N., Hayashi-Takagi, A., Ellis-Davies, G.C.R., Matsuzaki, M., Kasai, H., 2013. GABA promotes the competitive selection of dendritic spines by controlling local Ca2+ signaling. Nat. Neurosci. 16, 1409–1416.

He, K., Huertas, M., Hong, S.Z., Tie, X., Hell, J.W., Shouval, H., Kirkwood, A., 2015. Distinct eligibility traces for LTP and LTD in cortical synapses. Neuron 88, 528–538.

Hebb, D.O., 1949. The Organization of Behavior: A Neuropsychological Theory. Wiley, New York, NY, USA.

Helmholtz, H., 1925. Treatise on Physiological Optics, Vol. 3. Optical Society of America, Washington, DC.

Hodgkin, A.L., Huxley, A.F., 1952. A quantitative description of membrane current and its application to conduction and excitation in nerve. J. Physiol. 117 (4), 500–544.

Isomura, T., 2018. A measure of information available for inference. Entropy 20 (7), 512.

Isomura, T., Friston, K.J., 2018. In vitro neural networks minimise variational free energy. Sci. Rep. 8, 16926.

Isomura, T., Friston, K.J., 2020. Reverse-engineering neural networks to characterize their cost functions. Neural Comput. 32, 2085–2121.

Isomura, T., Kotani, K., Jimbo, Y., 2015. Cultured cortical neurons can perform blind source separation according to the free-energy principle. PLoS Comput. Biol. 11, e1004643.

Isomura, T., Parr, T., Friston, K.J., 2019. Bayesian filtering with multiple internal models: toward a theory of social intelligence. Neural Comput. 31 (12), 2390–2431.

Isomura, T., Shimazaki, H., Friston, K.J., 2022. Canonical neural networks perform active inference. Commun. Biol. 5, 55.

Johansen, J.P., Diaz-Mataix, L., Hamanaka, H., Ozawa, T., Ycu, E., Koivumaa, J., Kumar, A., Hou, M., Deisseroth, K., Boyden, E.S., LeDoux, J.E., 2014. Hebbian and neuromodulatory mechanisms interact to trigger associative memory formation. Proc. Natl. Acad. Sci. U. S. A. 111, E5584–E5592.

Kaplan, R., Friston, K.J., 2018. Planning and navigation as active inference. Biol. Cybern. 112, 323–343.

Kiebel, S.J., Daunizeau, J., Friston, K.J., 2008. A hierarchy of time-scales and the brain. PLoS Comput. Biol. 4, e1000209.

Knill, D.C., Pouget, A., 2004. The Bayesian brain: the role of uncertainty in neural coding and computation. Trends Neurosci. 27 (12), 712–719.

Kuśmierz, Ł., Isomura, T., Toyoizumi, T., 2017. Learning with three factors: modulating Hebbian plasticity with errors. Curr. Opin. Neurobiol. 46, 170–177.

Laje, R., Buonomano, D.V., 2013. Robust timing and motor patterns by taming chaos in recurrent neural networks. Nat. Neurosci. 16, 925–933.

Linsker, R., 1988. Self-organization in a perceptual network. Computer 21, 105–117.

London, M., Roth, A., Beeren, L., Häusser, M., Latham, P.E., 2010. Sensitivity to perturbations in vivo implies high noise and suggests rate coding in cortex. Nature 466 (7302), 123–127.

Maisto, D., Donnarumma, F., Pezzulo, G., 2015. Divide et impera: subgoaling reduces the complexity of probabilistic inference and problem solving. J. R. Soc. Interface 12, 20141335.

Malenka, R.C., Bear, M.F., 2004. LTP and LTD: an embarrassment of riches. Neuron 44, 5–21.

Markram, H., Lübke, J., Frotscher, M., Sakmann, B., 1997. Regulation of synaptic efficacy by coincidence of postsynaptic APs and EPSPs. Science 275, 213–215.

Millidge, B., 2020. Deep active inference as variational policy gradients. J. Math. Psychol. 96, 102348.

Newsome, W.T., Britten, K.H., Movshon, J.A., 1989. Neuronal correlates of a perceptual decision. Nature 341, 52–54.

Olshausen, B.A., Field, D.J., 1996. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. Nature 381 (6583), 607–609.

Paille, V., Fino, E., Du, K., Morera-Herreras, T., Perez, S., Kotaleski, J.H., Venance, L., 2013. GABAergic circuits control spike-timing-dependent plasticity. J. Neurosci. 33 (22), 9353–9363.

Parr, T., Friston, K.J., 2017. Uncertainty, epistemics and active inference. J. R. Soc. Interface 14 (136), 20170376.

Parr, T., Da Costa, L., Friston, K.J., 2020. Markov blankets, information geometry and stochastic thermodynamics. Phil. Trans. R. Soc. A 378, 20190159.

Pawlak, V., Wickens, J.R., Kirkwood, A., Kerr, J.N., 2010. Timing is not everything: neuromodulation opens the STDP gate. Front. Syn. Neurosci. 2, 146.

Pellicano, E., Burr, D., 2012. When the world becomes 'too real': a Bayesian explanation of autistic perception. Trends Cogn. Sci. 16 (10), 504–510.

Rao, R.P., Ballard, D.H., 1999. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. Nat. Neurosci. 2, 79–87.

Reynolds, J.N.J., Hyland, B.I., Wickens, J.R., 2001. A cellular mechanism of reward-related learning. Nature 413, 67–70.

Salgado, H., Köhr, G., Treviño, M., 2012. Noradrenergic 'tone' determines dichotomous control of cortical spike-timing-dependent plasticity. Sci. Rep. 2, 417.

Seol, G.H., Ziburkus, J., Huang, S., Song, L., Kim, I.T., Takamiya, K., Huganir, R.L., Lee, H.K., Kirkwood, A., 2007. Neuromodulators control the polarity of spike-timing-dependent synaptic plasticity. Neuron 55, 919–929.

Sussillo, D., Abbott, L.F., 2009. Generating coherent patterns of activity from chaotic neural networks. Neuron 63, 544–557.

Torigoe, M., Islam, T., Kakinuma, H., Fung, C.C.A., Isomura, T., Shimazaki, H., Aoki, T., Fukai, T., Okamoto, H., 2021. Zebrafish capable of generating future state prediction error show improved active avoidance behavior in virtual reality. Nat. Commun. 12, 5712.

Wald, A., 1947. An essentially complete class of admissible decision functions. Ann. Math. Stat. 18, 549–555.

Wieland, S., Schindler, S., Huber, C., Köhr, G., Oswald, M.J., Kelsch, W., 2015. Phasic dopamine modifies sensory-driven output of striatal neurons through synaptic plasticity. J. Neurosci. 35, 9946–9956.

Wolpert, D.M., Kawato, M., 1998. Multiple paired forward and inverse models for motor control. Neural Netw. 11 (7–8), 1317–1329.

Yagishita, S., Hayashi-Takagi, A., Ellis-Davies, G.C.R., Urakubo, H., Ishii, S., Kasai, H., 2014. A critical time window for dopamine actions on the structural plasticity of dendritic spines. Science 345, 1616–1620.

Zhang, J.C., Lau, P.M., Bi, G.Q., 2009. Gain in sensitivity and loss in temporal contrast of STDP by dopaminergic modulation at hippocampal synapses. Proc. Natl. Acad. Sci. U. S. A. 106, 13028–13033.