

脳知能の統一理論に基づく AGI開発の展望



磯村 拓哉

理化学研究所脳神経科学研究センター 脳型知能理論研究ユニット
第9回全脳アーキテクチャ・シンポジウム 2024年9月18日

磯村 拓哉

理化学研究所 脳神経科学研究センター 脳型知能理論研究ユニット ユニットリーダー

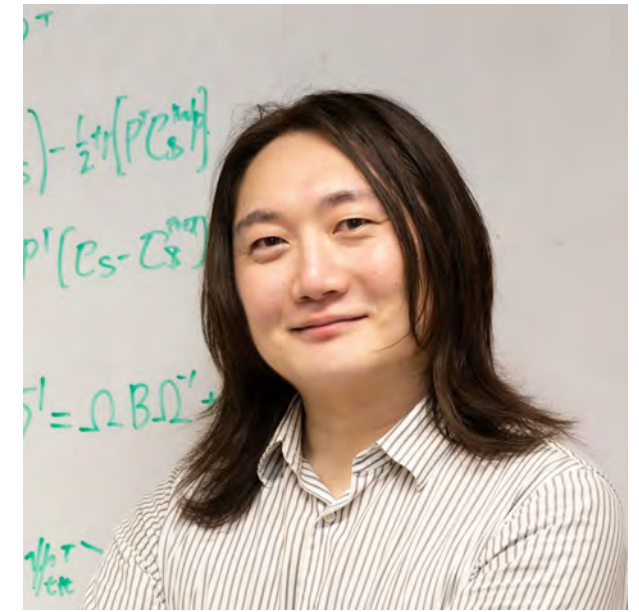
京都大学大学院情報学研究科 連携准教授

文科省科研費 学術変革領域 (A) 「予測と行動の統一理論の開拓と検証 (2023~2027)」 領域代表

takuya.isomura@riken.jp

@TakuyaIsomura

<https://cbs.riken.jp/jp/faculty/t.isomura/>



専攻：理論神経科学

研究テーマ：脳の知能が持つ普遍的な特性を数学を使って表現すること

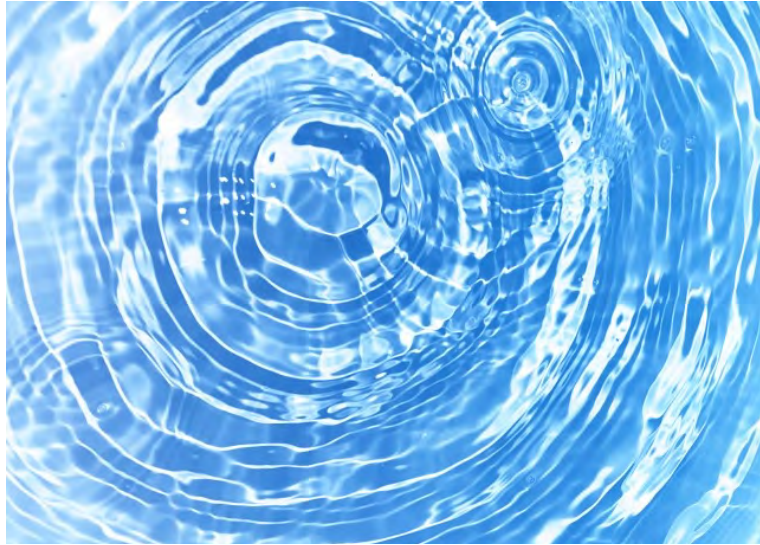
研究の動機

- 生物のような人工知能を作りたい
- そのために生物の知能を理解したい

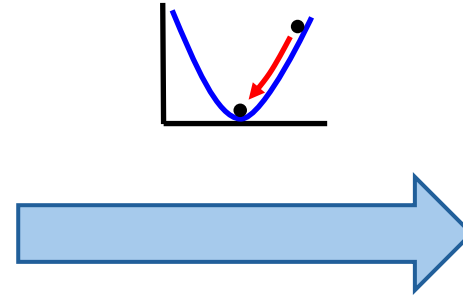
問い

- 生物の知能の本質的原理は？
- 生物の脳が機械より優れている点は？

物理学



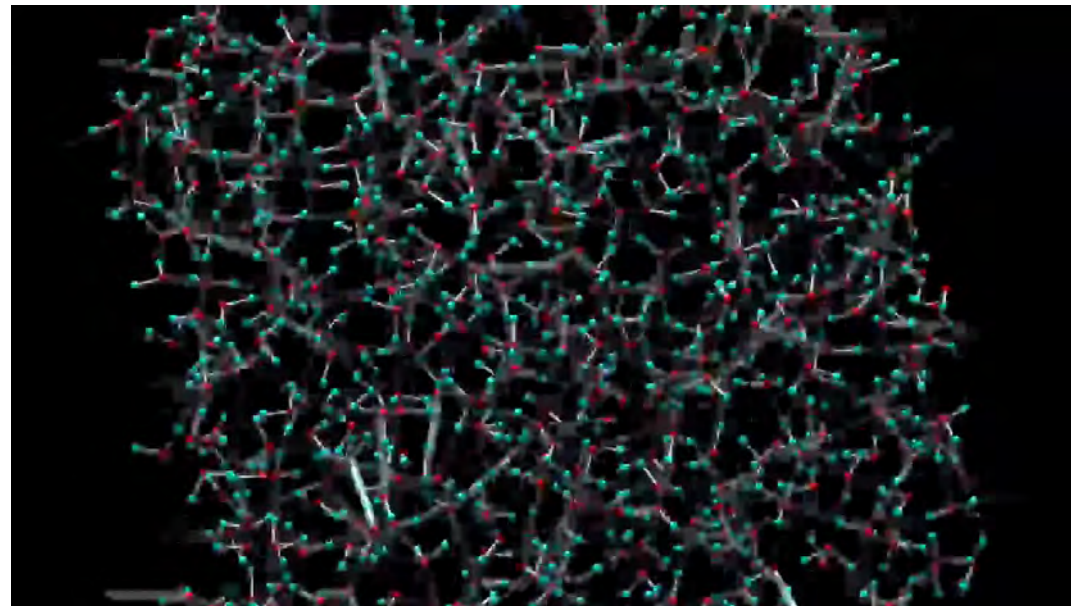
無秩序・無規則



エネルギーを下げる



規則性・パターンが生まれる

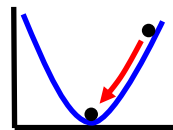


分子科学研究所, 水が氷になるまで, <https://www.youtube.com/watch?v=8eXdXHP5dk8>

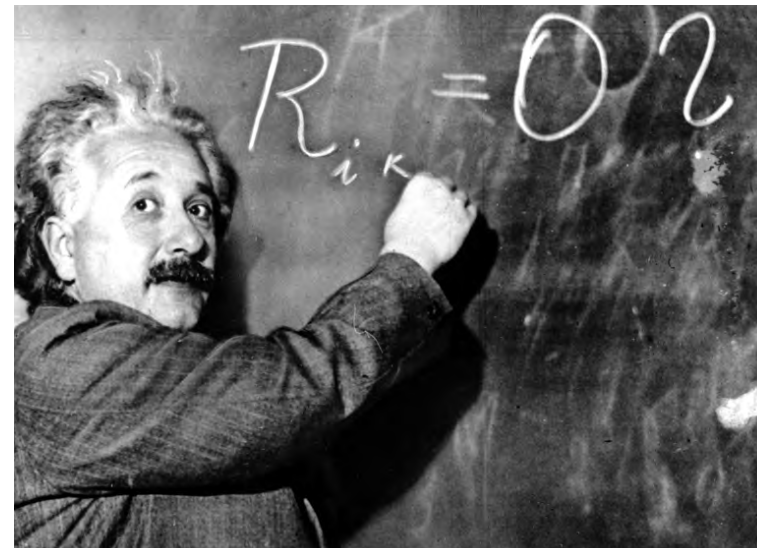
知能の科学



赤ちゃん (学習前)



エネルギーを下げる



知能が生まれる

生まれたばかりの脳の神経回路



外界の情報をほとんど持っていない



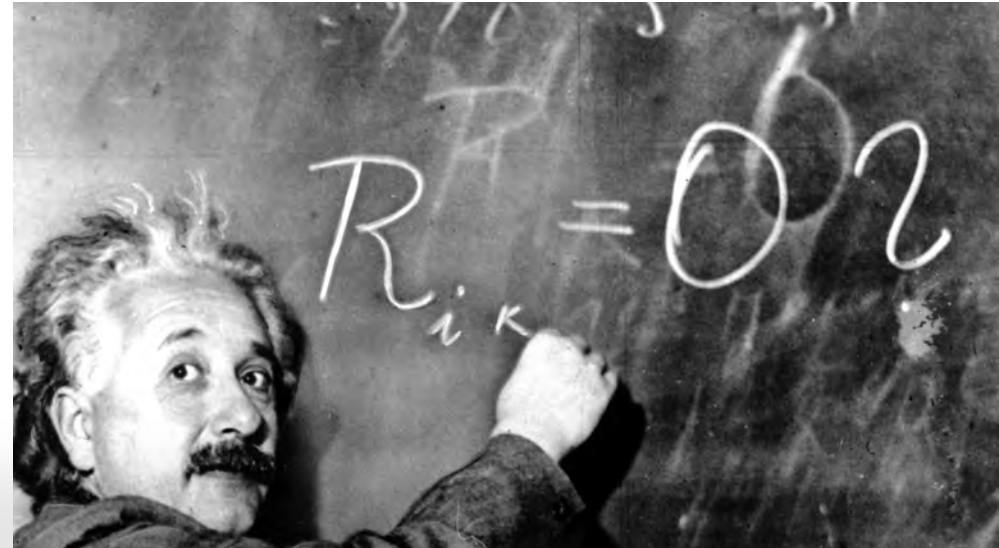
学習

外界に適応した脳の神経回路



外界に応じたパターンが生まれる
予測・洞察・創造ができるようになる

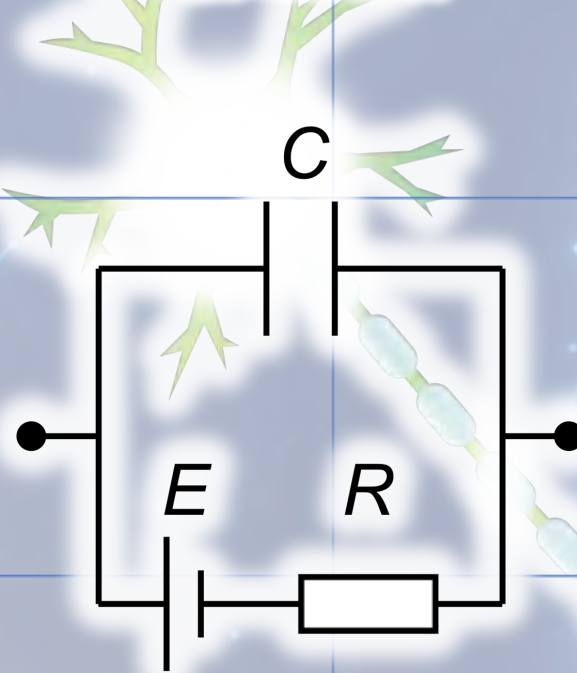
知能とはなにか？



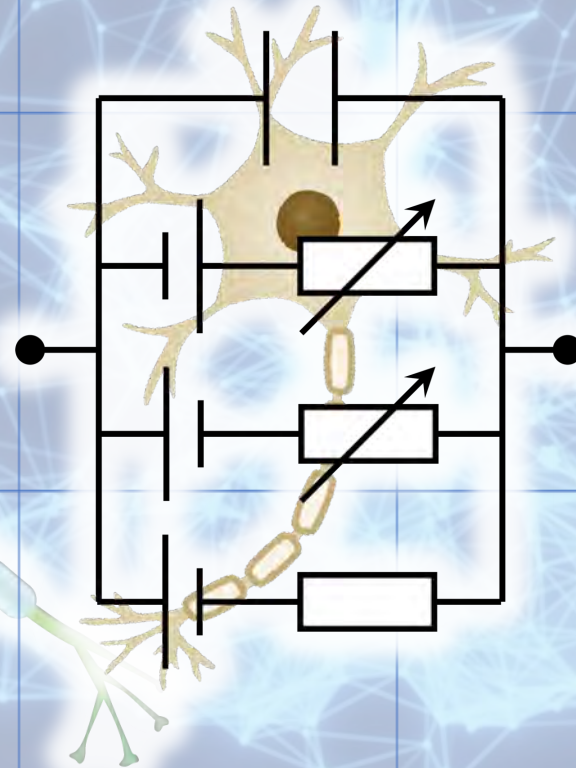
脳の統一理論



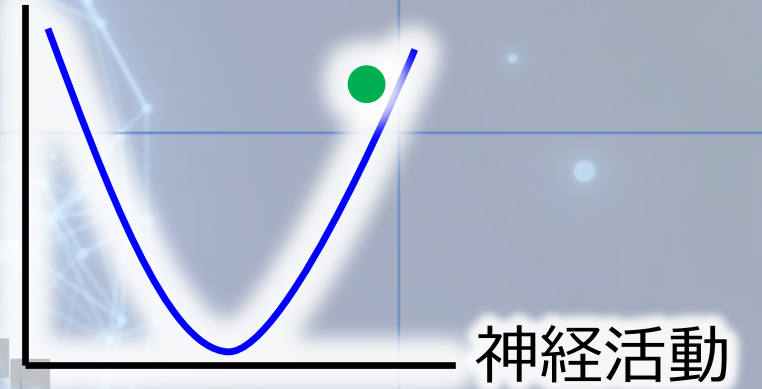
電気回路



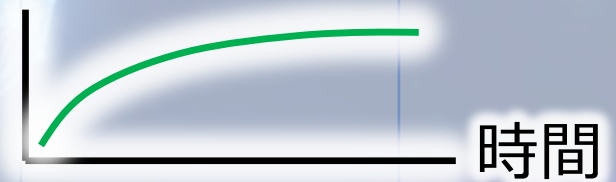
神経細胞



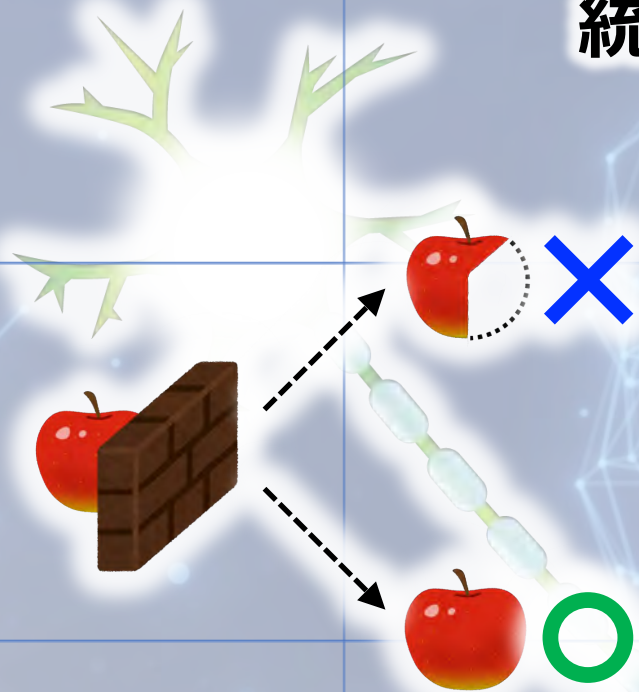
神経回路のコスト関数



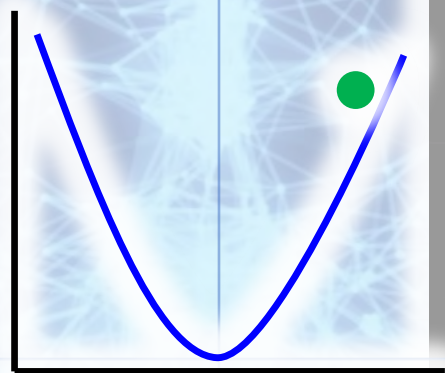
神経活動



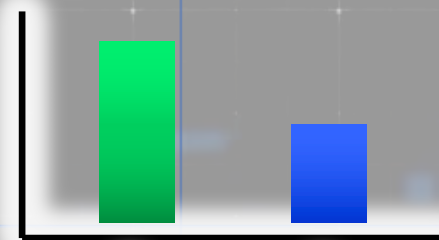
統計学的な推論 (ベイズ推論)



自由エネルギー



期待値



期待値



$$P(\text{state}|\text{input}) = \frac{P(\text{input}|\text{state})P(\text{state})}{P(\text{input})}$$

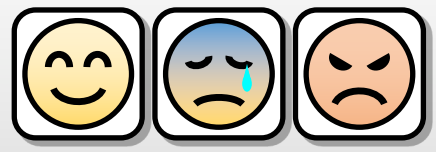
自由エネルギー原理

- Karl J. Fristonが提唱している脳の情報理論
- 生物の知覚や学習、行動は、変分自由エネルギーと呼ばれるコスト関数を最小化するように決まるとしている
- その結果、生物は変分ベイズ推論と呼ばれる統計学的な推論を自己組織化的に行うとされている

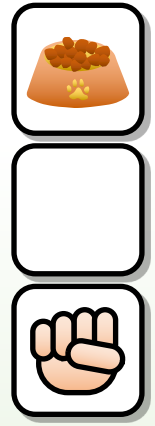
Friston, Nat Rev Neurosci, 2010

パラメータ θ

隠れ状態 s



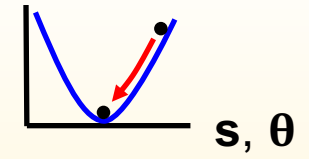
感覚入力 o



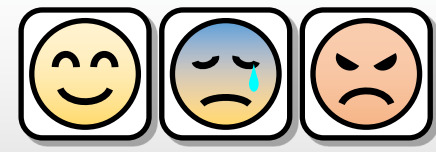
$$F = E_Q[-\ln P(o, s, \theta) + \ln Q(s, \theta)]$$

生成モデル

自由エネルギー F



事後期待値 s

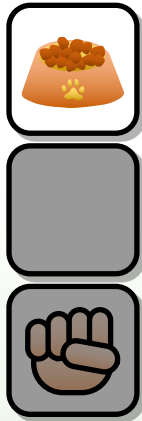


パラメータ θ

隠れ状態 s

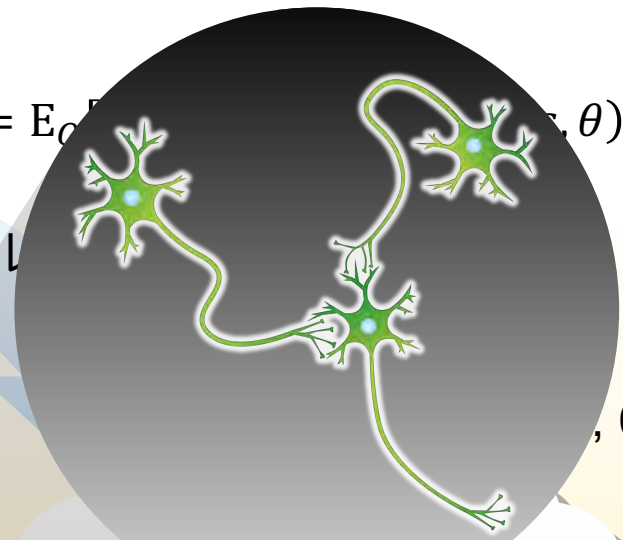


感覚入力 o

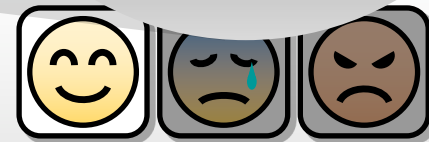


生成モデル

$$F = E_{\phi} [F(o, s, \theta)]$$



θ



好みの事前分布 C

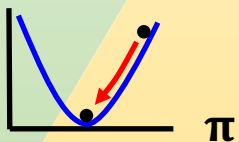


行動 δ

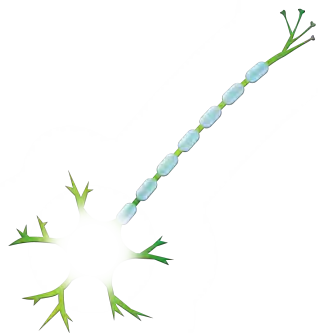
行動方策 π



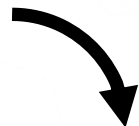
期待自由エネルギー G



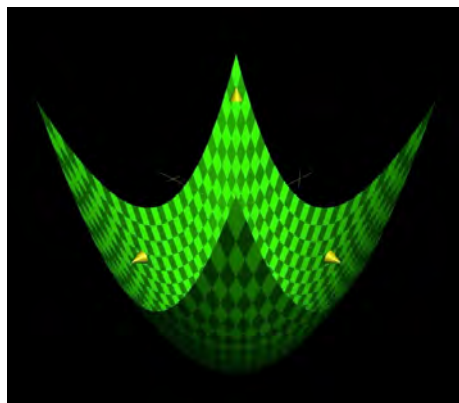
神経活動の方程式



積分

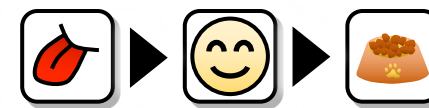


ヘルムホルツエネルギー \mathcal{A}

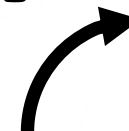


脳が予測や学習を行う方法

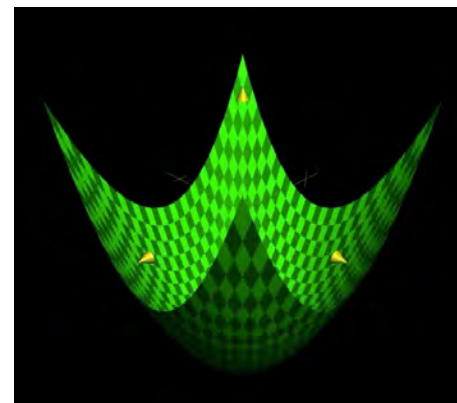
ベイズ推論



最小化



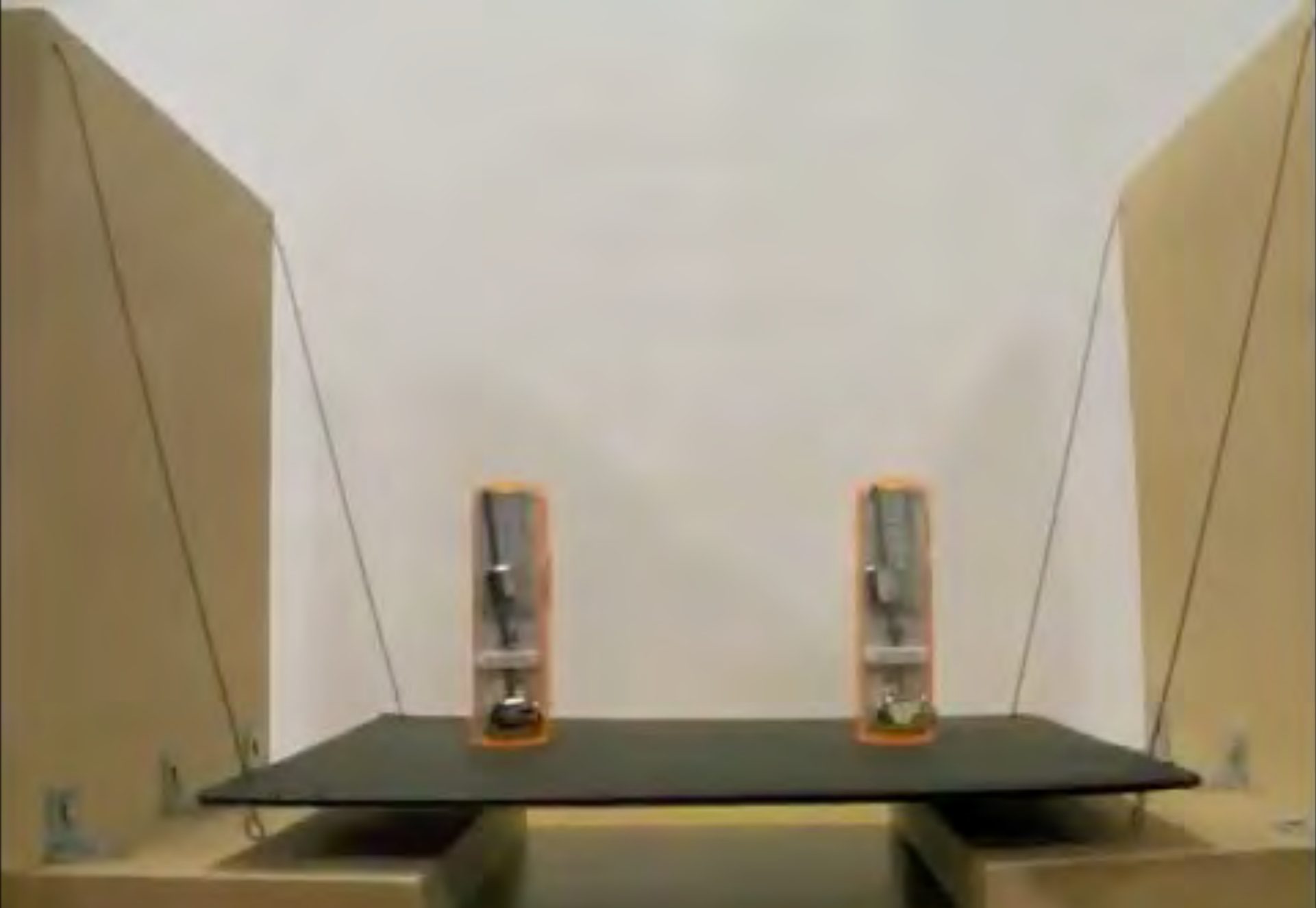
変分自由エネルギー \mathcal{F}



AIが予測や学習を行う方法と同じ

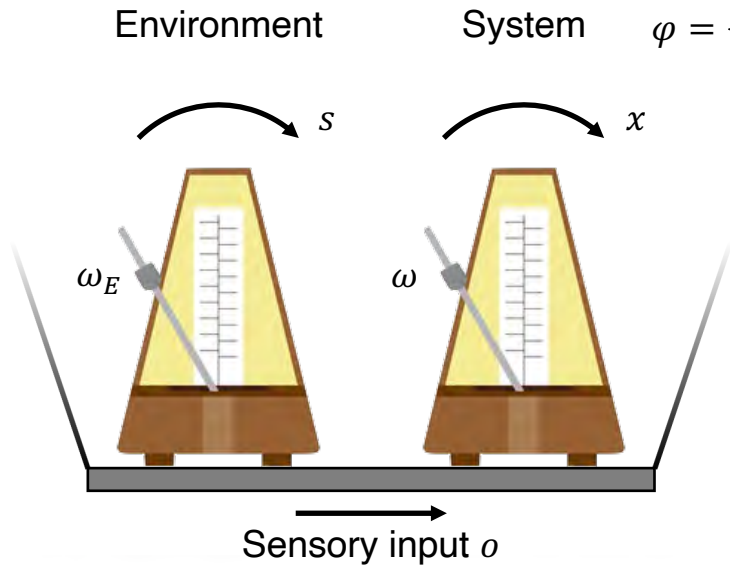


数学的に等価



Synchronization of two metronomes, Filmed at Ikeguchi Laboratory
<https://www.youtube.com/watch?v=feEBzjqishQ>

自己組織化系のベイズ力学：どんな力学系もベイズ推論と見なせる



状態の経路とパラメータのダイナミクスは虚数時間勾配法に従う

$$\partial_\tau \varphi = -\partial_\varphi H + \xi$$

対応するFokker-Planck方程式は次のHelmholtz energyを最小化している

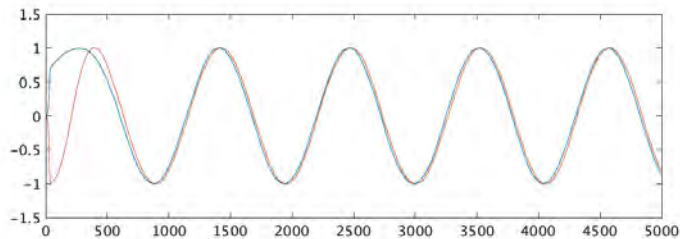
$$\mathcal{A}[\pi(\varphi), \tilde{\sigma}] = \left\langle H(\varphi, \tilde{\sigma}) + \frac{1}{\beta} \log \pi(\varphi) \right\rangle_{\pi(\varphi)}$$

変分自由エネルギー $\mathcal{F}[q(\vartheta), \tilde{\sigma}] = \langle -\log p_m(\vartheta, \tilde{\sigma}) + \log q(\vartheta) \rangle_{q(\vartheta)}$

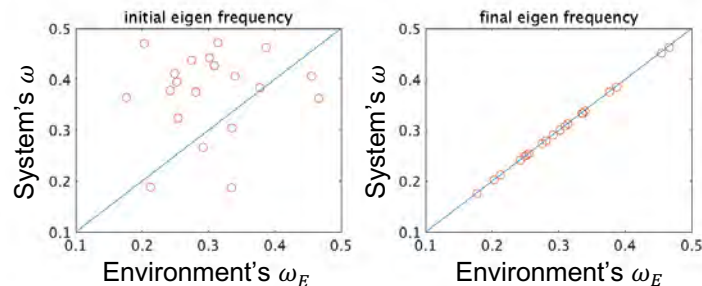
$\beta = 1$ のとき、 \mathcal{A} と \mathcal{F} は数式として等しい（自然同値）（c.f., 完備類定理）

\mathcal{F} を最小化すると事後分布が得られる（ベイズの定理） $q(\vartheta) = \frac{p_m(\vartheta, \tilde{\sigma})}{p_m(\tilde{\sigma})}$

力学系の経路とパラメータが共通の \mathcal{A} を最小化することを要請するだけで計算機構が自己組織化(進化)的に創発する可能性を示唆



Matching of eigenfrequency



スケール or
抽象度

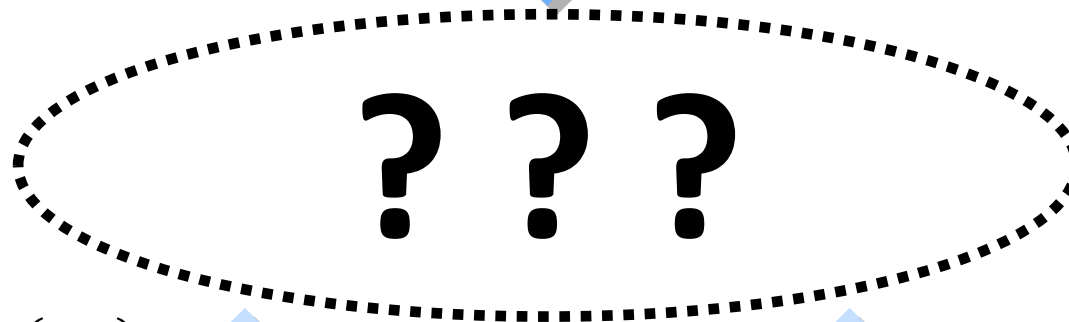


H. von Helmholtz

自由エネルギー原理
a.k.a. ベイズカ学



K. J. Friston



$$\dot{x} = \text{func}(x, o)$$

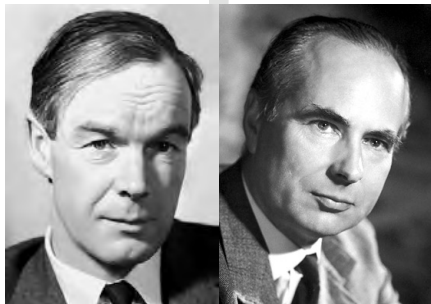
発火率モデル

ヘップ可塑性

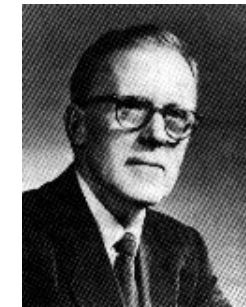
$$\dot{W} = \text{pre} \times \text{post}$$

Hodgkin-Huxley eq.

Spike-timing
dependent plasticity

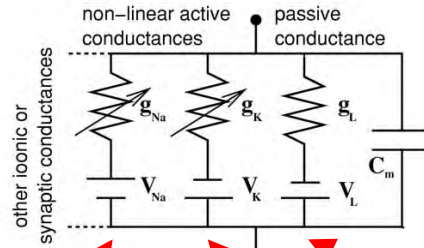
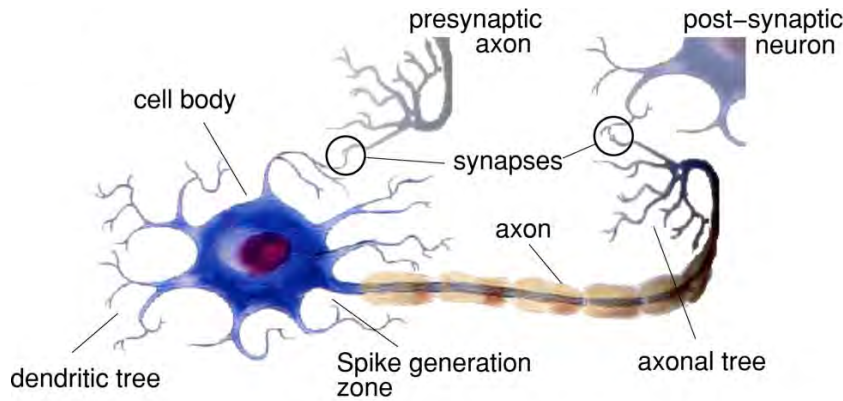


A. L. Hodgkin & A. F. Huxley

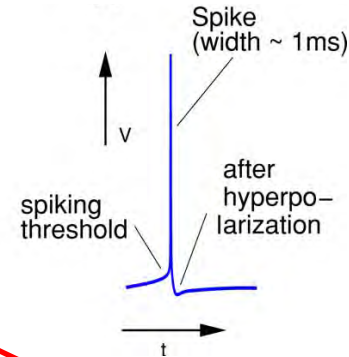


D. O. Hebb

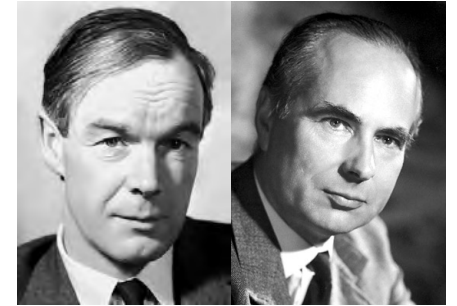
神経細胞の活動とシナプス可塑性



Rabinovich, Varona, Selverston & Abarbanel, *Rev Mod Phys*, 2006



1963年ノーベル賞

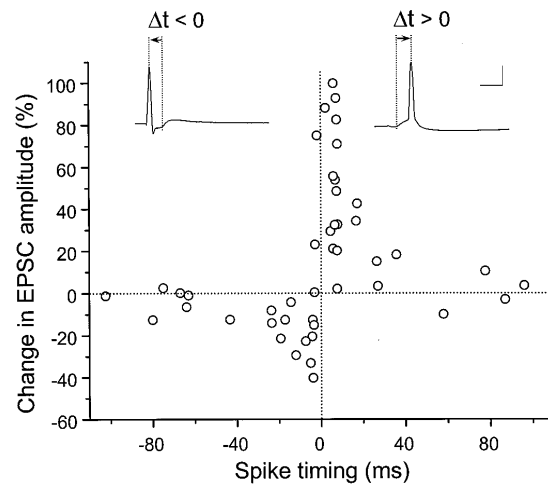


Hodgkin & Huxley, 1952

Hodgkin-Huxley方程式

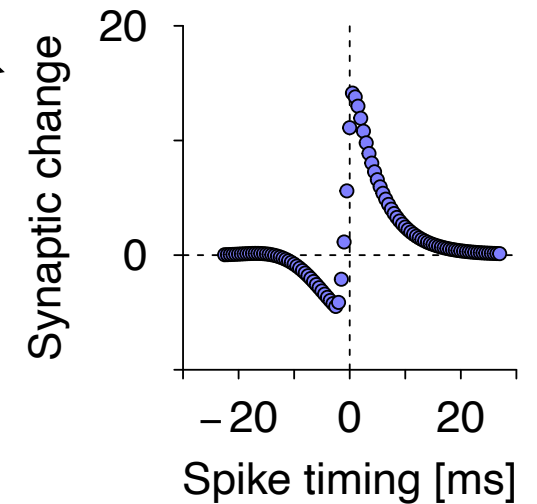
$$\begin{cases} C\dot{v} = -g_{Na}m^3h(v - V_{Na}) - g_Kn^4(v - V_K) - g_L(v - V_L) + I \\ \dot{m} = (m_\infty(v) - m)/\tau_m(v) \\ \dot{h} = (h_\infty(v) - h)/\tau_h(v) \\ \dot{n} = (n_\infty(v) - n)/\tau_n(v) \end{cases}$$

スパイクタイミング依存
シナプス可塑性



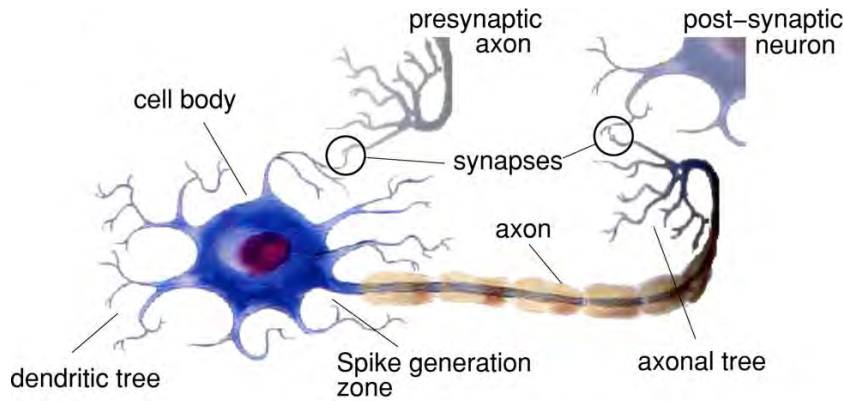
Bi & Poo, *J Neurosci*, 1998

ベイズ力学
による理論予想
Isomura, *arXiv*, 2023

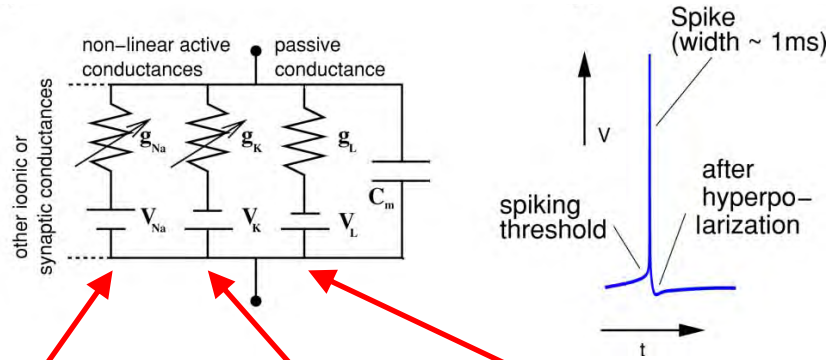


同じ形

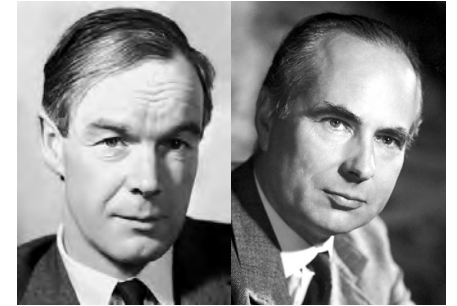
神経細胞の活動とシナプス可塑性



Rabinovich, Varona, Selverston & Abarbanel, Rev Mod Phys, 2006



1963年ノーベル賞



Hodgkin & Huxley, 1952

Hodgkin-Huxley方程式

縮約

$$\begin{cases} C\dot{v} = -g_{Na}m^3h(v - V_{Na}) - g_Kn^4(v - V_K) - g_L(v - V_L) + I \\ \dot{m} = (m_\infty(v) - m)/\tau_m(v) \\ \dot{h} = (h_\infty(v) - h)/\tau_h(v) \\ \dot{n} = (n_\infty(v) - n)/\tau_n(v) \end{cases}$$

2D Hodgkin-Huxley方程式

$$\begin{cases} \dot{v} \propto f_v(v, u) + I \\ \dot{u} \propto f_u(v, u) \end{cases}$$

フィッツフュー・南雲モデル
Canonical neuron model

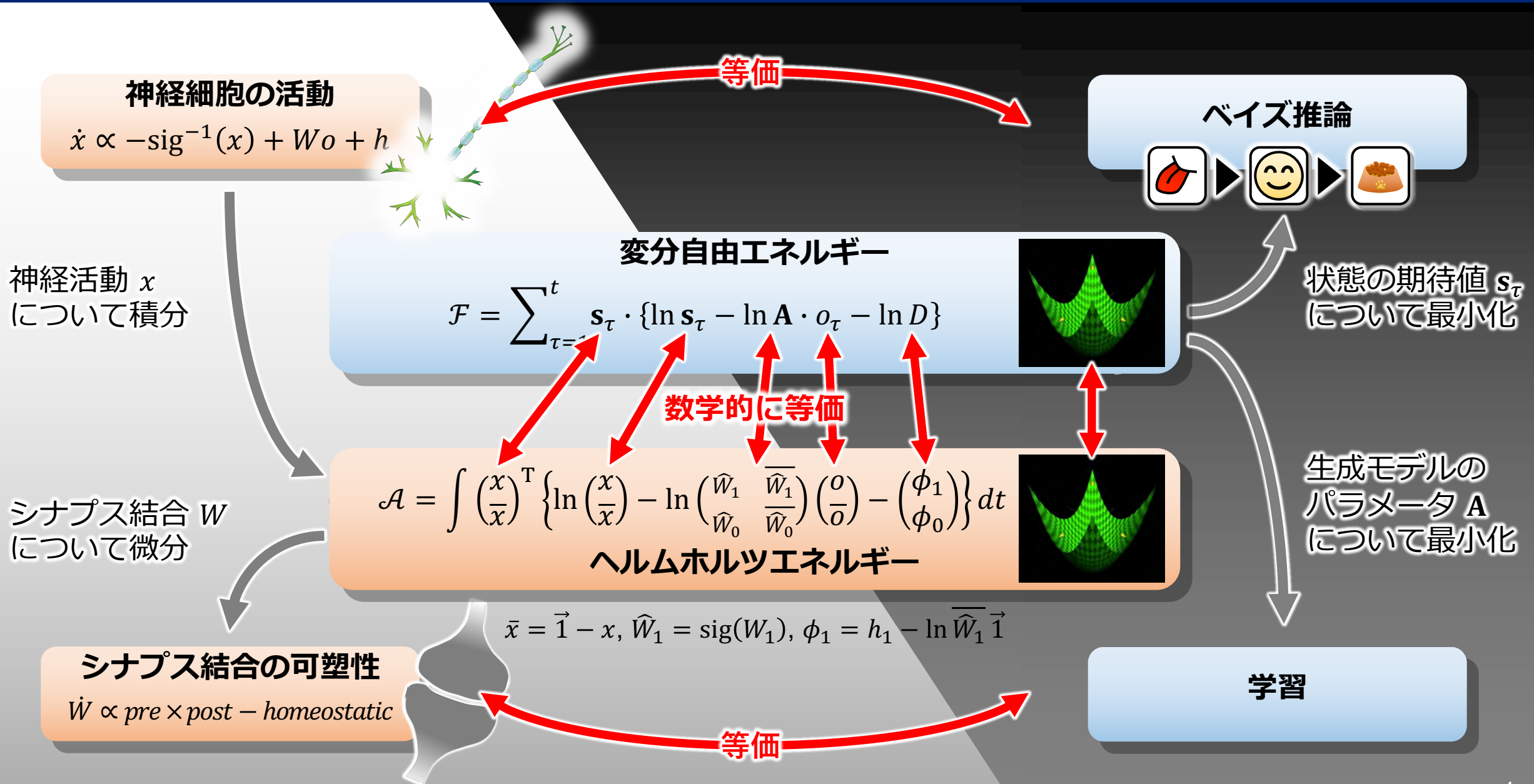
結合

正準神経回路モデル (canonical neural networks)

$$\dot{x}_t = \underbrace{-\text{sig}^{-1}(x_t)}_{\text{leak current}} + \underbrace{(W_1 - W_0)o_t}_{\text{synaptic input}} + \underbrace{h_1 - h_0}_{\text{threshold}}$$

x_t : 神経活動, o_t : 感覚入力, W : シナプス結合強度, h : 発火閾値

“理論上は”どんな神経回路もベイズ推論を行っている



生成モデルのリバーズエンジニアリングを用いた検証

実物脳（神経回路の活動）

人工脳（ベイズ推論）

生成モデルのリバーズ
エンジニアリング

1 神経細胞の活動を計測

回路構造の同定

2 神経活動モデルの割り当て

積分

定数の推定

3 神経回路のコスト関数の同定

数理的な等価性

4 生成モデル&自由エネルギーの同定

微分

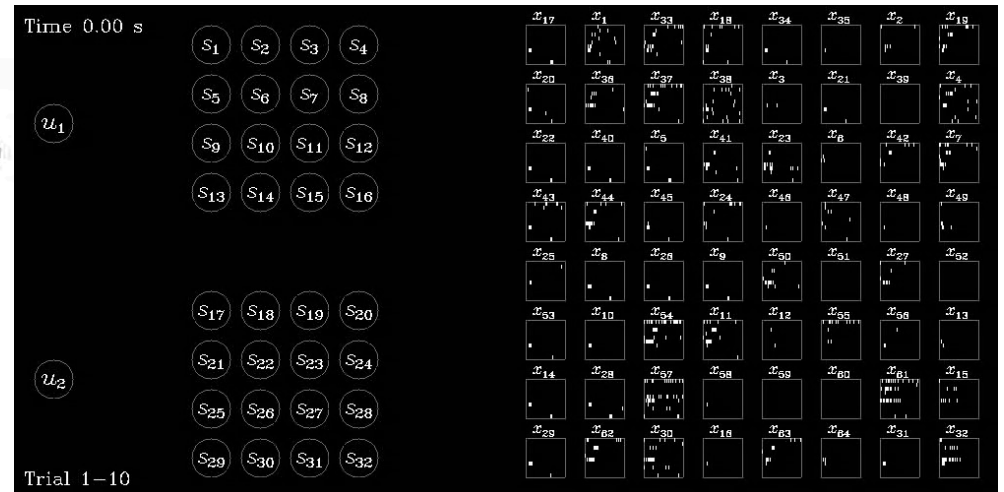
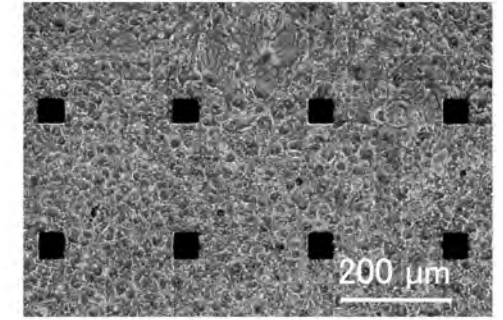
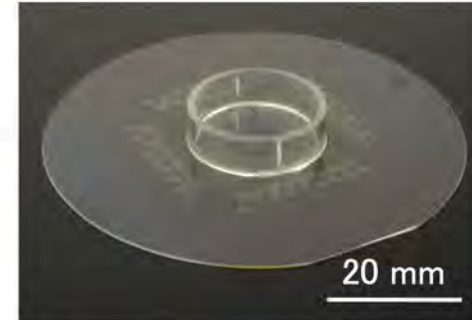
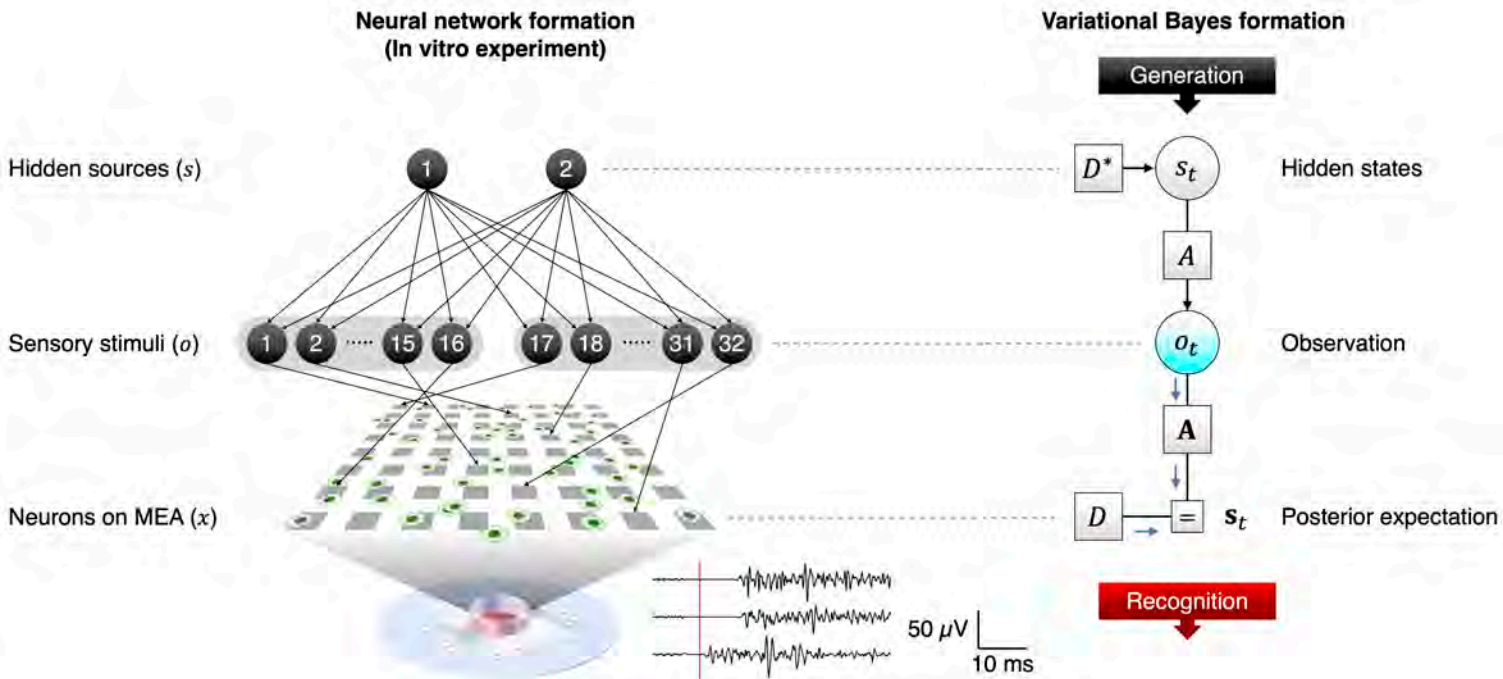
5 学習アルゴリズムの導出

時間積分

6 学習結果の予測

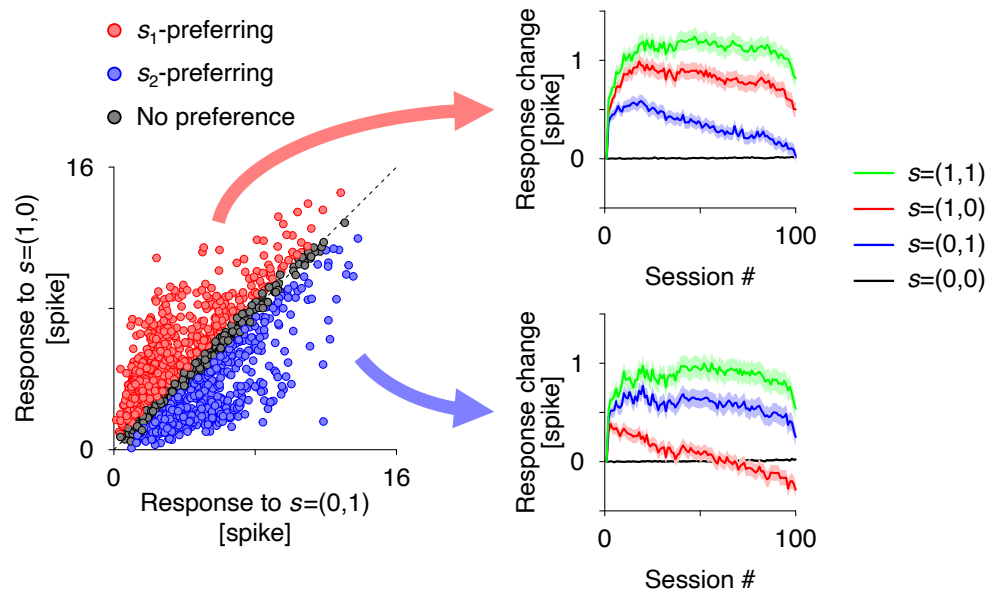
実験結果の予測(診断)

培養神経回路における自由エネルギー原理の実証

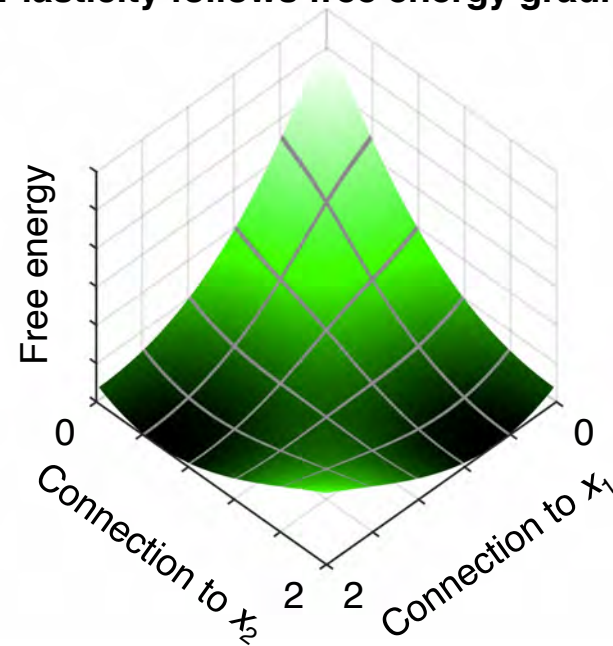


培養神経回路における自由エネルギー原理の実証

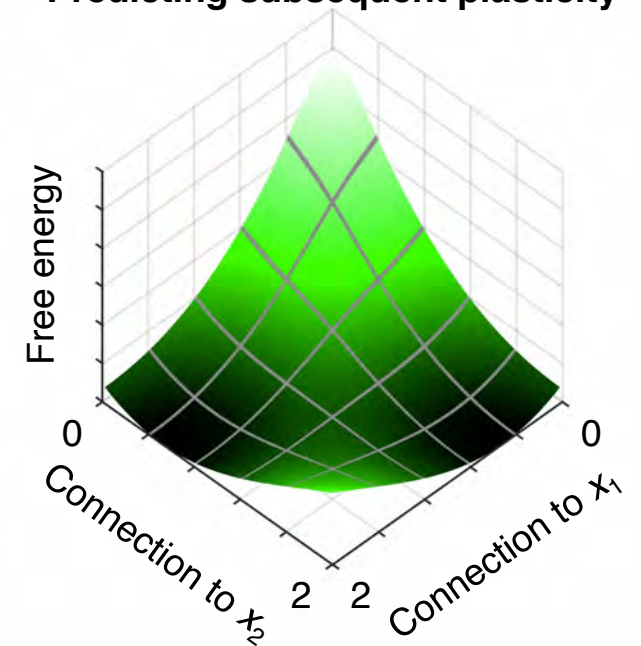
Response preference to hidden sources



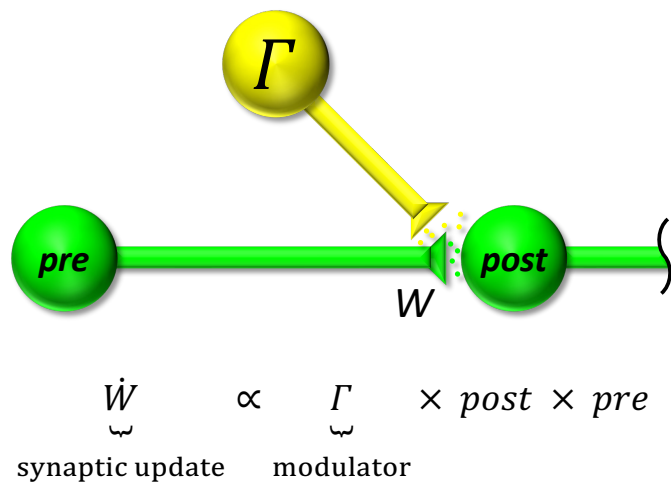
Plasticity follows free energy gradient



Predicting subsequent plasticity

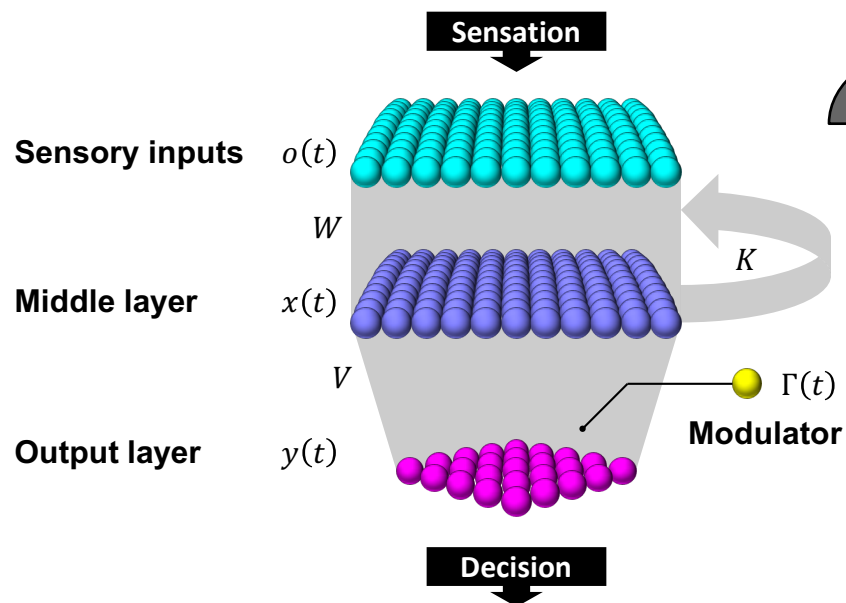
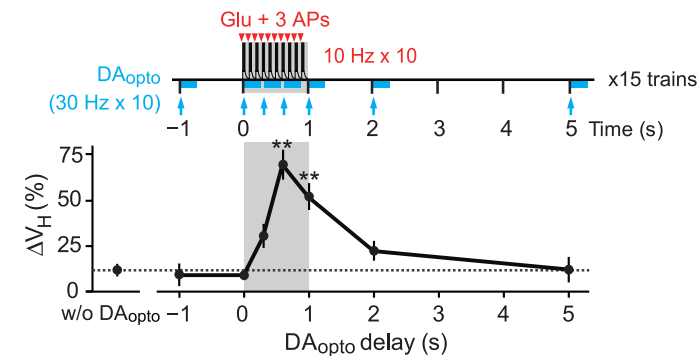


Canonical neural networks perform active inference



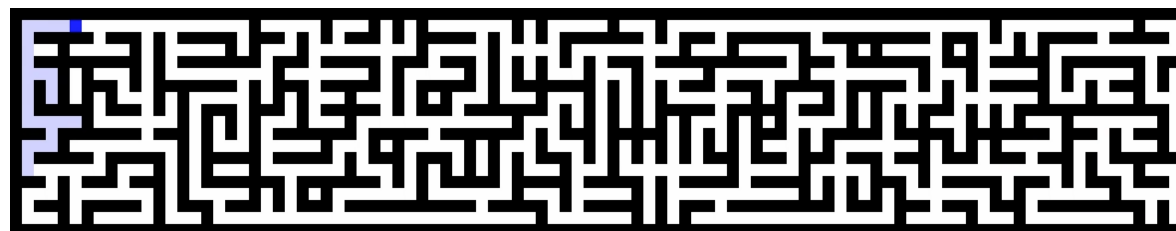
Delayed modulation of plasticity

Yagishita et al., Science, 2014



Delayed modulation corresponds to post hoc evaluation of past decisions

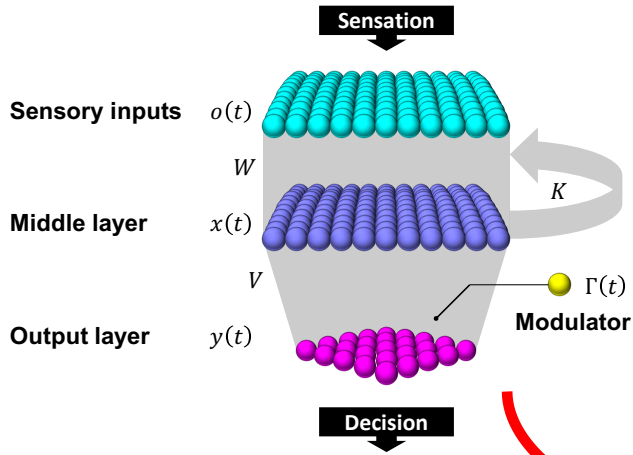
Trajectory of the agent in a maze



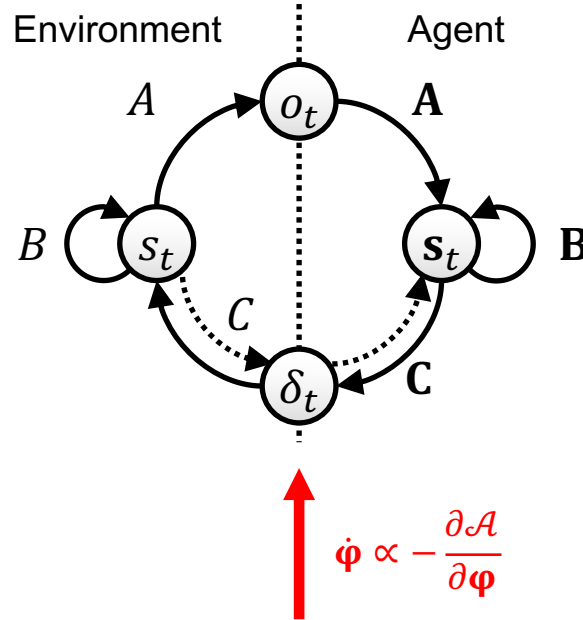
Isomura et al., Commun Biol, 2022

Equivalence between canonical NNs, Bayesian inference, and Turing machines

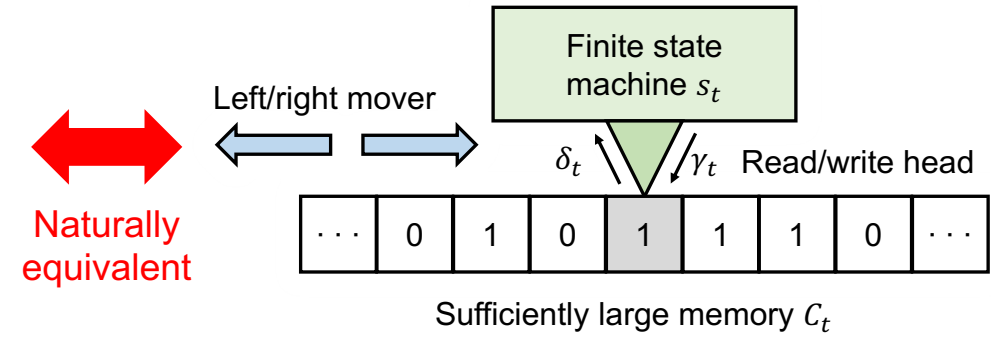
Canonical neural networks



Variational Bayes under POMDP



Differentiable Turing machine



Integral

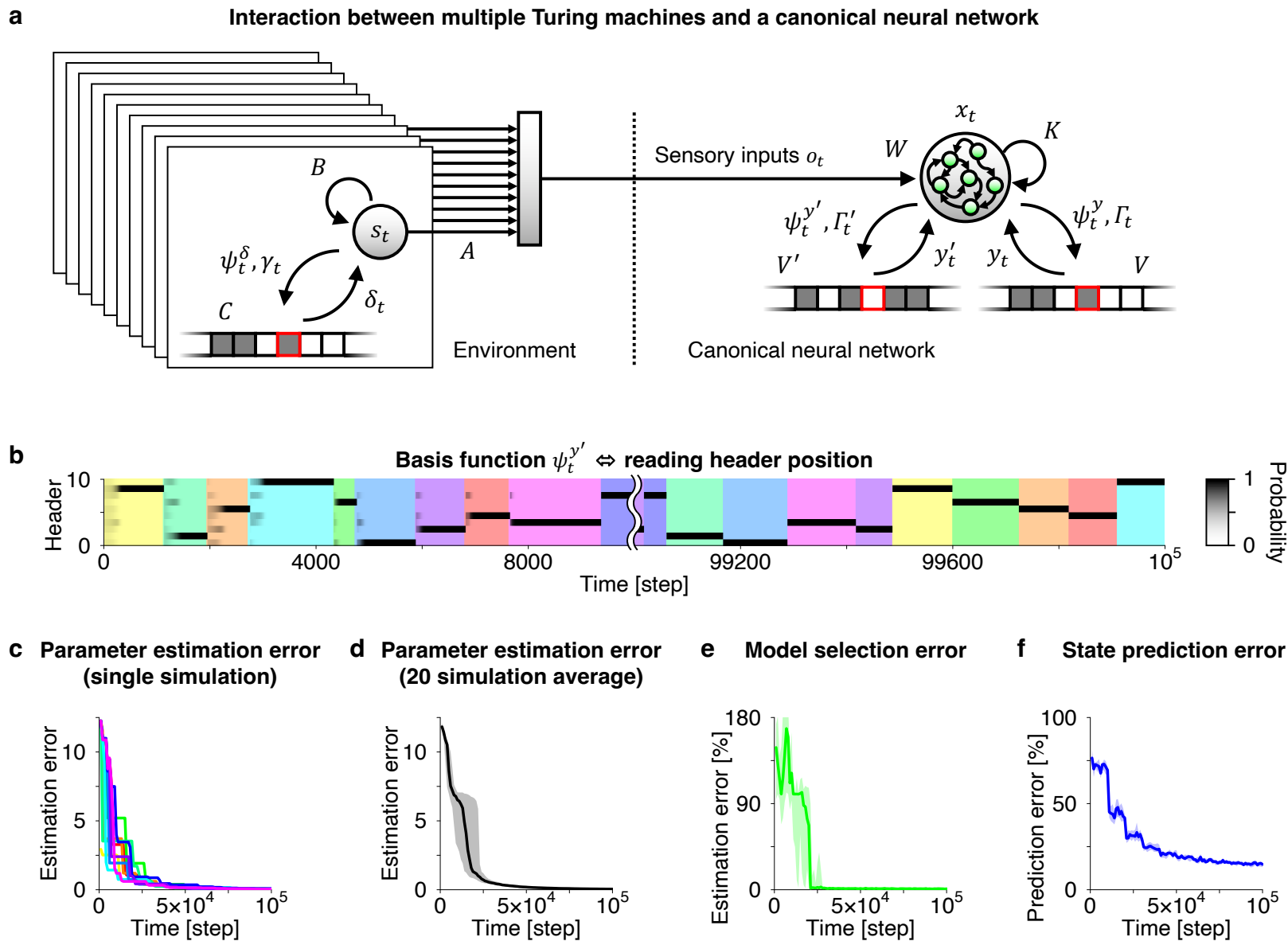
$$\dot{\phi} \propto -\frac{\partial \mathcal{A}}{\partial \phi}$$

$$\dot{\phi} \propto -\frac{\partial \mathcal{A}}{\partial \phi}$$

Helmholtz energy for canonical NNs

$$\mathcal{A}[\pi(\phi), o_{1:t}, \xi] = \sum_{\tau=1}^t \left(\frac{x_{\tau}}{\bar{x}_{\tau}} \right) \cdot \left\{ \ln \left(\frac{x_{\tau}}{\bar{x}_{\tau}} \right) - \left(\frac{W_1}{W_0} \right) \cdot o_{\tau} - \left(\frac{K_1}{K_0} \right) \psi_{\tau-1}^x - \begin{pmatrix} h_1^x \\ h_0^x \end{pmatrix} \right\} + \sum_{\tau=1}^t \left(\frac{y_{\tau}}{\bar{y}_{\tau}} \right) \cdot \left\{ \ln \left(\frac{y_{\tau}}{\bar{y}_{\tau}} \right) - (\vec{1} - 2\Gamma_{\tau}) \odot \begin{pmatrix} V_1 \\ V_0 \end{pmatrix} \psi_{\tau-1}^y - \begin{pmatrix} h_1^y \\ h_0^y \end{pmatrix} \right\} + c$$

Canonical neural networks can implement universal Turing machines



まとめ

- 自由エネルギー原理：全ての生物の知覚や学習、行動は、変分自由エネルギーを最小化するように決まり、その結果ベイズ推論を自己組織化的に行うという主張。
- 等価性：正準神経回路のダイナミクスは、変分自由エネルギーの最小化をしていると見なすことができる。外界のベイズ推論を行うことは、神経回路の普遍的な特性。
- 実証実験：いくつかの実験において、神経回路の自己組織化を理論的に定量予測可能。
- 完全性：正準神経回路は神経活動と可塑性により万能チューリングマシンを実装可能。

展望：脳型学習アルゴリズムのAI応用

現在のAIの学習アルゴリズム（バックプロパゲーション）は多くの訓練データと計算量が必要（データ枯渇問題）

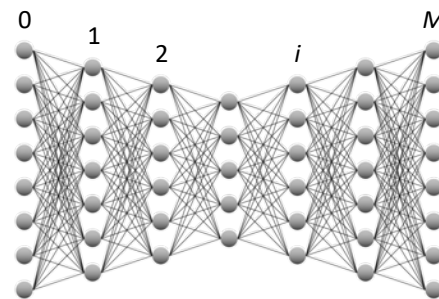


脳のような効率良い学習アルゴリズム

高いデータ効率、ローカルな計算のみ、ノイズが大きくても可、逐次的に学習可、悪い解に陥ることを回避

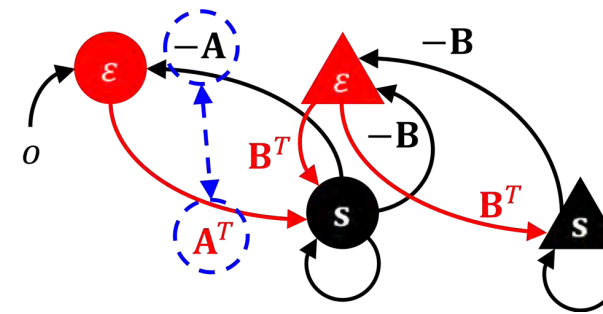
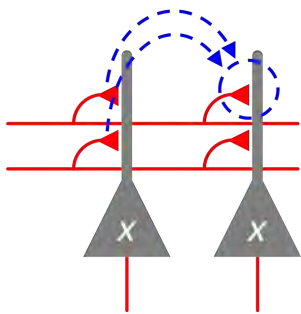
バックプロパゲーション

$$\dot{W}_i \propto - \left\langle \frac{\partial l}{\partial W_i} \right\rangle = - \left\langle \frac{\partial u_M^T}{\partial u_i} \frac{\partial l}{\partial u_M} f_{i-1}^T \right\rangle \\ = \langle f_i' W_{i+1}^T \cdots f_{M-1}' W_M^T (x_0 - x_M) f_{i-1}^T \rangle$$



予測符号化

$$\dot{s} - Ds \propto - \frac{\partial \mathcal{F}}{\partial s} \\ \dot{A} \propto - \frac{\partial \bar{\mathcal{F}}}{\partial A}, \quad \dot{B} \propto - \frac{\partial \bar{\mathcal{F}}}{\partial B}$$



別の脳型学習アルゴリズムの可能性