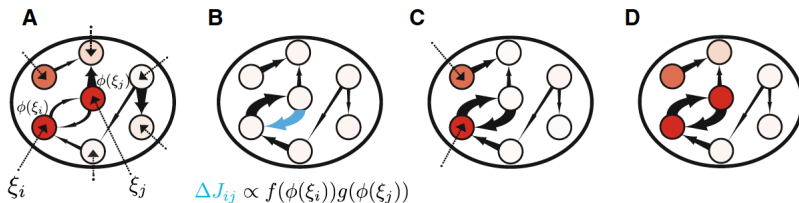# Information bounds and attractor dynamics of a Hebbian associative memory

Clayton Seitz

May 26, 2021
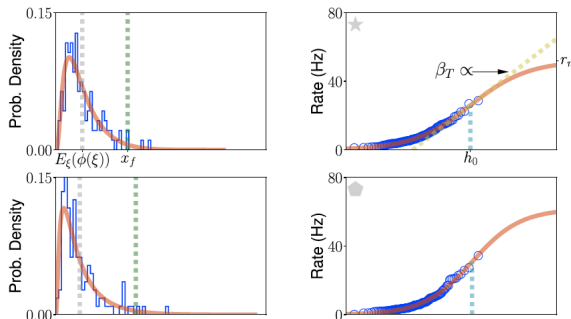
# RNNs trained with a Hebbian learning rule

It is difficult to infer learning rules *in vivo* solely from spikes



$$\Delta J_{ij} \propto f(\phi(\xi_i))g(\phi(\xi_j))$$

Learning rules can be inferred (with assumptions) from firing rates [1]

---

[1][Peirera and Brunel, Neuron. 2018]

Measuring the *static* transfer function from novel images assuming that input currents are Gaussian variables
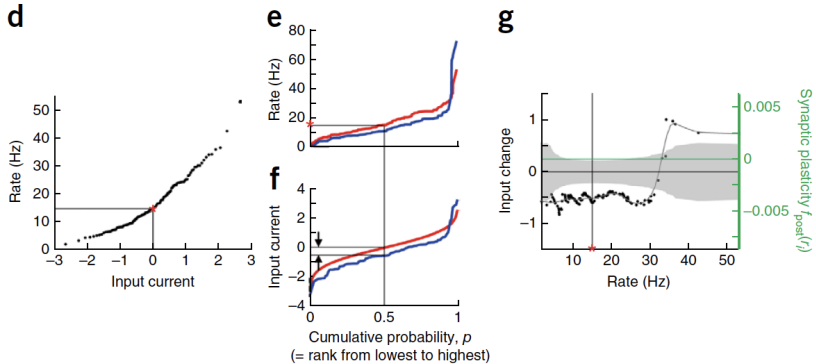
$$\phi(\boldsymbol{\xi}) = \frac{r_{max}}{1 + \exp \beta(\boldsymbol{\xi} - \boldsymbol{\xi}_0)}$$

[2]

[2][Peirera and Brunel, Neuron. 2018]

**d**

Rate (Hz): 0, 10, 20, 30, 40, 50

Input current: −3, −2, −1, 0, 1, 2, 3

**e**

Rate (Hz): 0, 20, 40, 60, 80

**f**

Input current: −4, −2, 0, 2, 4

Cumulative probability, $p$
(= rank from lowest to highest): 0, 0.5, 1

**g**

Input change: −1, 0, 1

Synaptic plasticity $f_{post}(r)$: −0.005, 0, 0.005

Rate (Hz): 10, 20, 30, 40, 50

Inferring the change in input current $\xi_{in}$ from the change in firing rate in novel relative to familiar stimuli

[3]

---

[3][Lim et al., Nature Neuroscience. 2015]

**d**

**e**

**f**

**g**

The change in input current to a neuron can then be read from the firing rate of that neuron when presented a novel stimulus
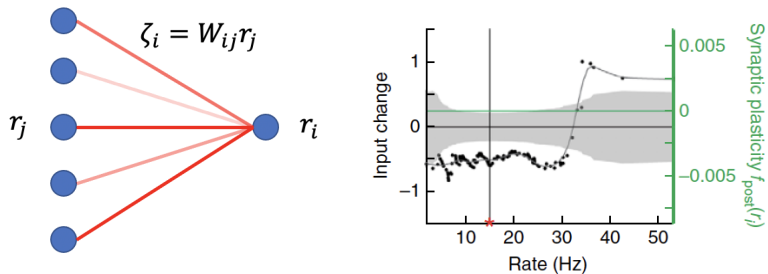
$$\Delta\xi_i(r) \propto (2q + 1 - \tanh(\beta(r - x)))$$

[4]

[4][Lim et al., Nature Neuroscience. 2015]

$$\zeta_i = W_{ij} r_j$$

Assuming that $\Delta W_{ij} \propto f(r_i) g(r_j)$, the change in input current $\xi_i$ is related to the learning rule by

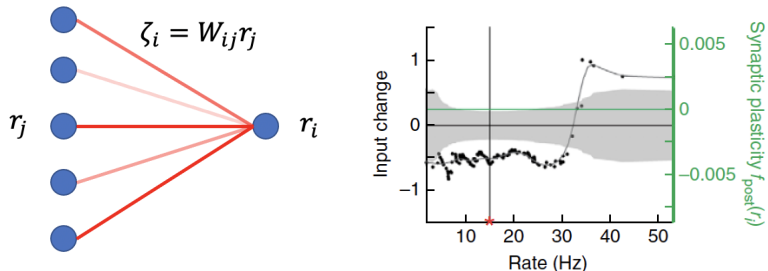$$\Delta \xi_i \propto f(r_i) \sum_j g(r_j) r_j$$

# Determining the learning rule



We can fit input change $\Delta\xi_i(r)$ from the data (top right)

$$\Delta\xi_i(r_i) \propto (2q + 1 - \tanh(\beta(r_i - r_0)))$$

5

---
[5][Lim et al., Nature Neuroscience. 2015]

# Determining the learning rule



We can infer the dependence of the learning rule on the post-synaptic firing rate $f(r_i)$

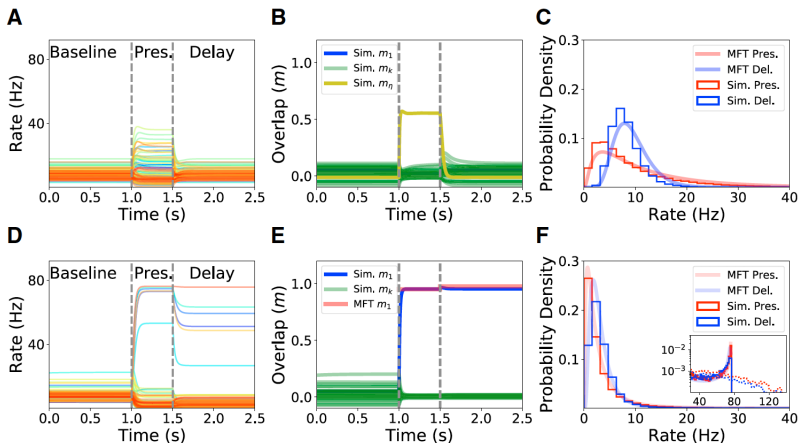$$f(r_i) = \Delta \xi_i(r_i) / \sum_j g(r_j) r_j$$

[6][Lim et al., Nature Neuroscience. 2015]

During training, we stimulate the network with

$$\xi_{in}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2}|\boldsymbol{\Sigma}|^{1/2}} \exp -\frac{1}{2}(\boldsymbol{r} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{r} - \boldsymbol{\mu})$$

# Attractor dynamics of the trained model



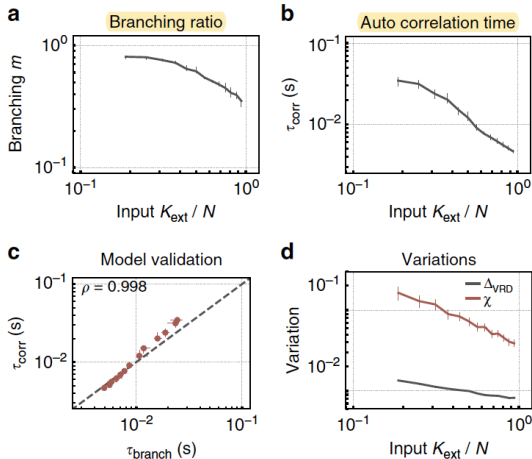Attractor states can be observed from the overlap $m$

[7][Peirera and Brunel, Neuron. 2018]

# Do these networks optimize information transmission?

Are these networks functioning at a critical point? What about the balance between input and recurrence? (Cramer et al. 2020)

How much information does the response $R$ carry about the input pattern $S$ i.e. $I(R; S)$ on novel and familiar stimuli?

What is the fundamental coding capacity of these networks?