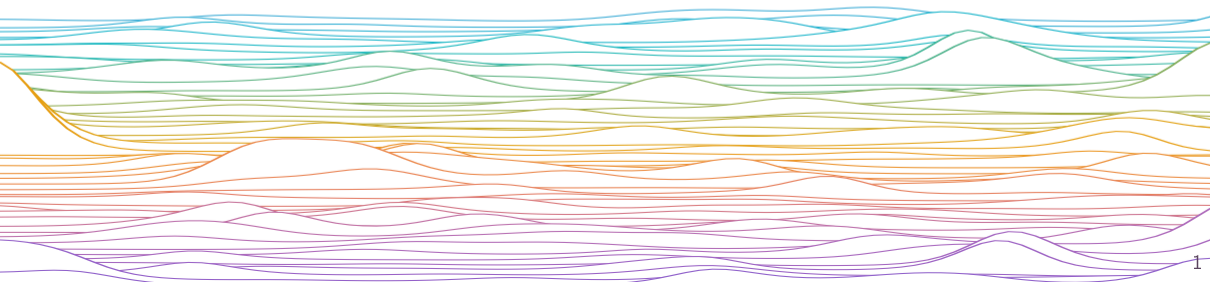


Self Healing Codes

5 November 2020

Michael Rule



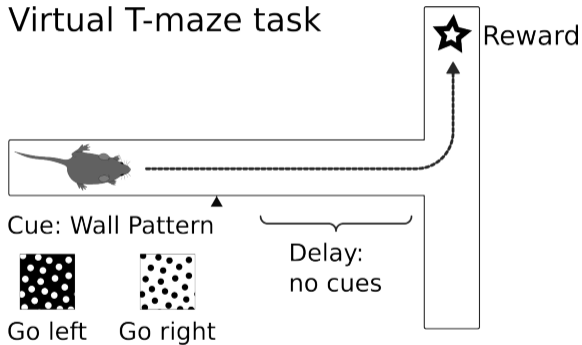
The brain is plastic.

The brain is stable.

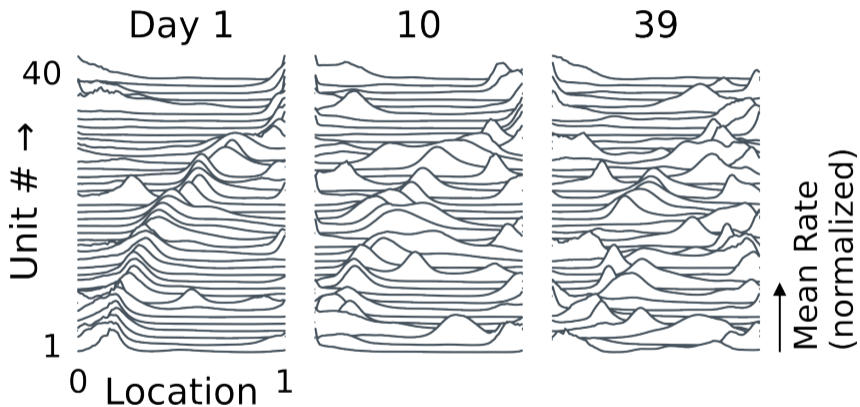
Image neural population codes over time



Virtual T-maze task

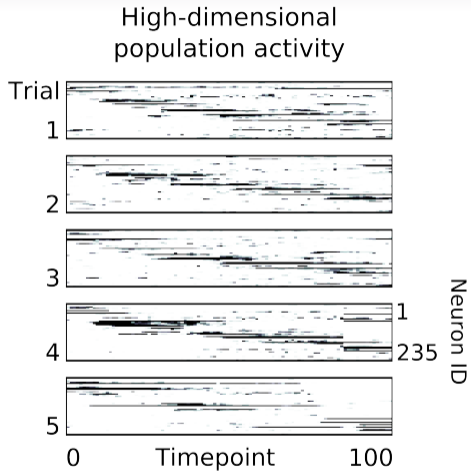


Neural population code is unstable*

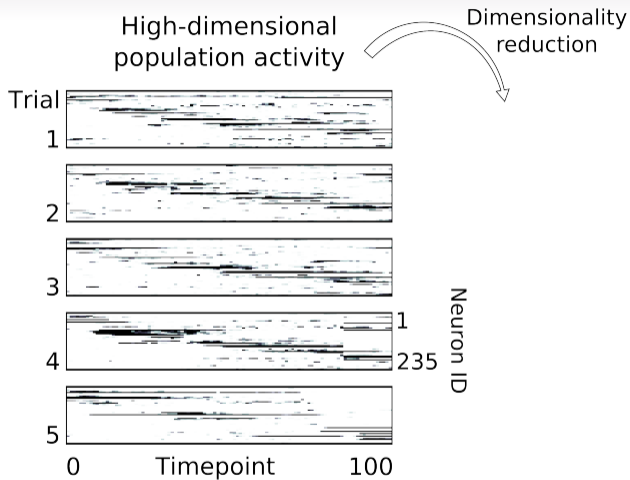


*In Posterior Parietal Cortex (PPC)

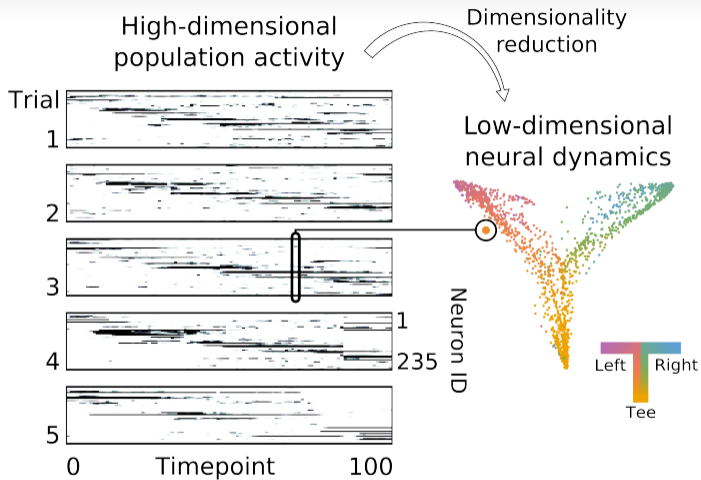
Representational Drift



Rule, O'Leary, Harvey,(2019) *Causes and consequences of representational drift.*



Rule, O'Leary, Harvey, (2019) *Causes and consequences of representational drift.*



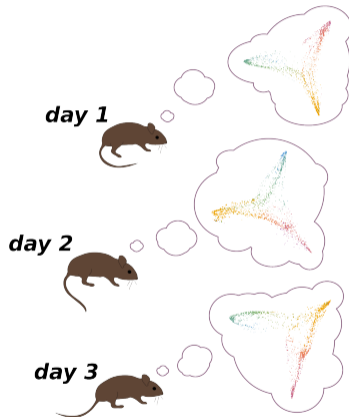
Rule, O'Leary, Harvey, (2019) *Causes and consequences of representational drift.*

Many degrees of freedom
in internal representations



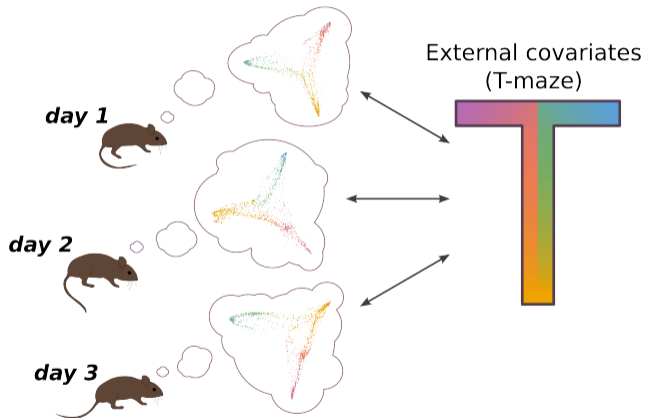
Rule, O'Leary, Harvey,(2019) *Causes and consequences of representational drift.*

Many degrees of freedom
in internal representations



Rule, O'Leary, Harvey, (2019) *Causes and consequences of representational drift.*

Many degrees of freedom
in internal representations



Rule, O'Leary, Harvey, (2019) *Causes and consequences of representational drift.*

Stable Task Information

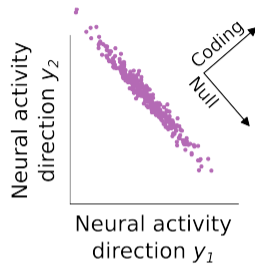
Rule ME, Loback AR, Raman DV, Driscoll L, Harvey CD, O'Leary T. 2020. Stable task information from an unstable neural population. eLife.

How is the brain is robust to changing (Δ) neural codes?

How is the brain is robust to changing (Δ) neural codes?

Invariance:

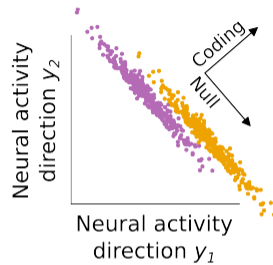
- Δ in null-space of readout



How is the brain is robust to changing (Δ) neural codes?

Invariance:

- Δ in null-space of readout



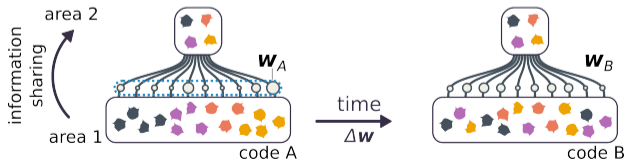
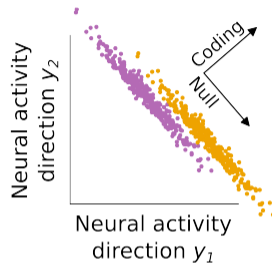
How is the brain is robust to changing (Δ) neural codes?

Invariance:

- Δ in null-space of readout

Coordination:

- Slow Δ , readout adapts



How is the brain is robust to changing (Δ) neural codes?

Invariance:

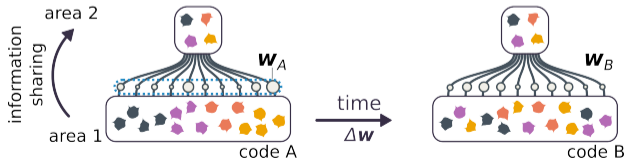
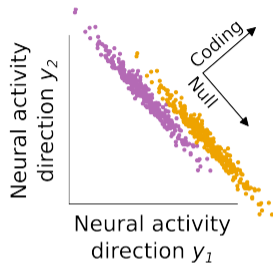
- Δ in null-space of readout

Coordination:

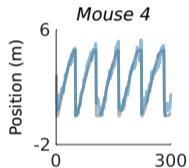
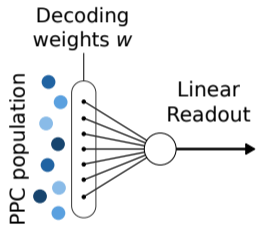
- Slow Δ , readout adapts

Analyse Driscoll et al. '17

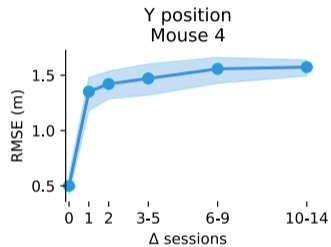
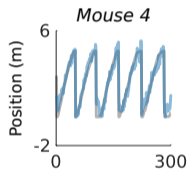
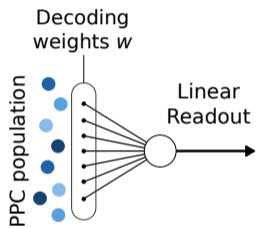
- Δ preserves invariant readout
- Slow plasticity can track Δ



Single-day decoders generalize poorly



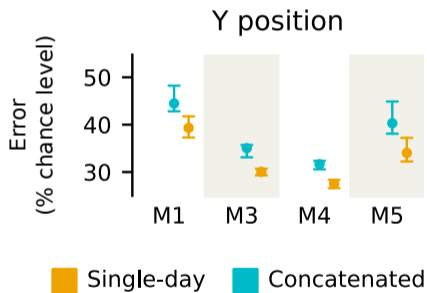
Single-day decoders generalize poorly



... but hint at long-term stable structure

Long-term \approx stable subspace exists, drift is constrained

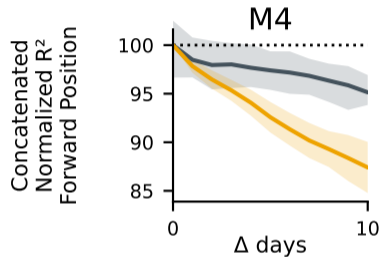
Fixed decoder trained over data from 7-10 days nearly as good as single-day



Long-term \approx stable subspace exists, drift is constrained

Fixed decoder trained over data from 7-10 days nearly as good as single-day

Unconstrained drift rapidly degrades performance



■ Data

■ Null model (random drift)

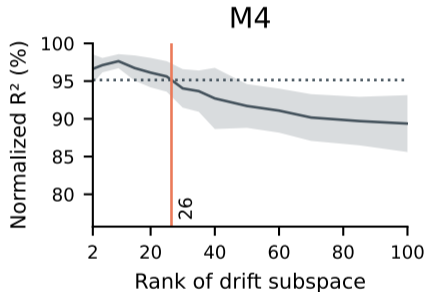
Shaded = 95% confidence

Long-term \approx stable subspace exists, drift is constrained

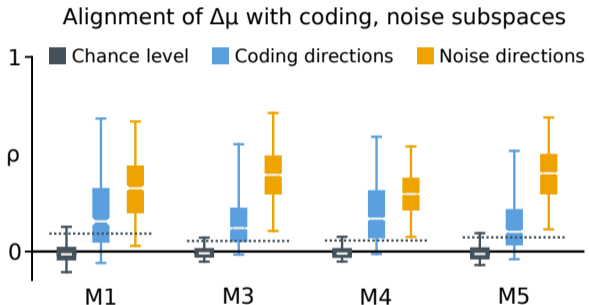
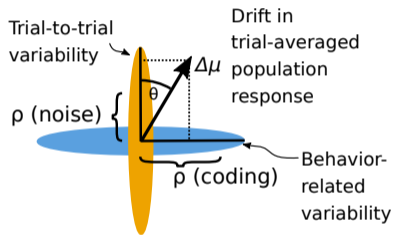
Fixed decoder trained over data from 7-10 days nearly as good as single-day

Unconstrained drift rapidly degrades performance

Results consistent with low-rank drift



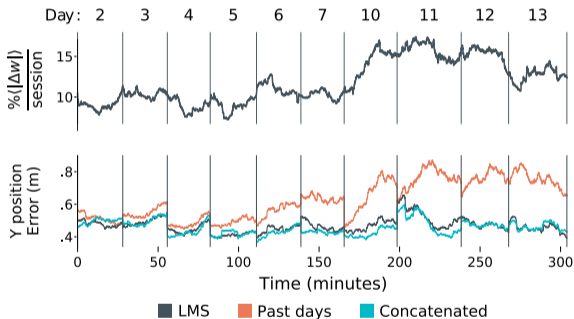
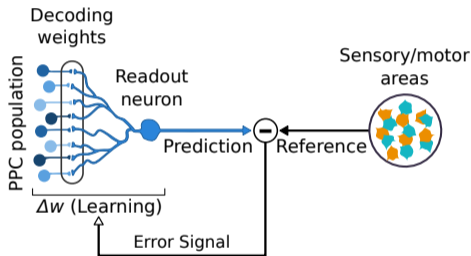
Drift resembles trial-to-trial variability



... But some drift occurs in directions that encode task information

≈ Stable subspace can be identified, tracked with modest plasticity

Distributed representations could detect tuning changes, adjust decoding weights



(~10-15% weight change per session for ~100 cells, more cells → less plasticity)

Drift (Δ) is structured

- Long-term Δ less than expected
- Consistent with low-rank
- More Δ in null directions

Track \approx stable subspace

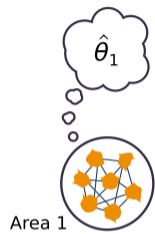
- Non-null Δ : slow and easy to track
- Weak error feedback sufficient

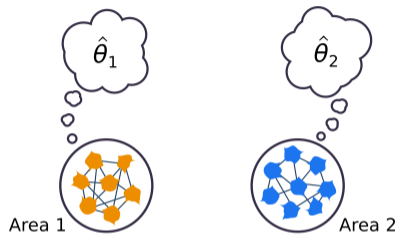
This talk: **Could neurons use what is stable to track what is volatile?**

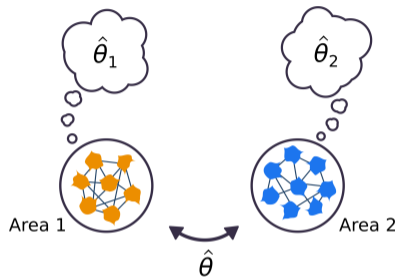
- Learning and correlated activation is sufficient to track drift
- How to stabilize readout without external feedback

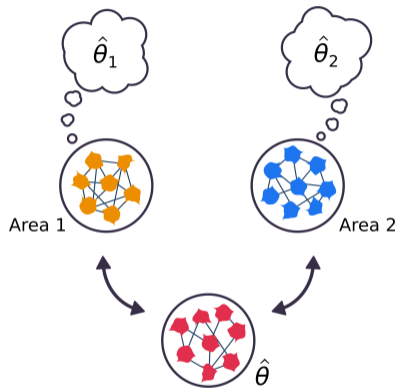
Ongoing learning addresses drift



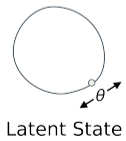


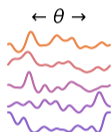
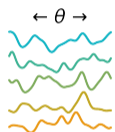






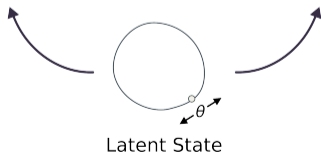
- Low-D latent variable



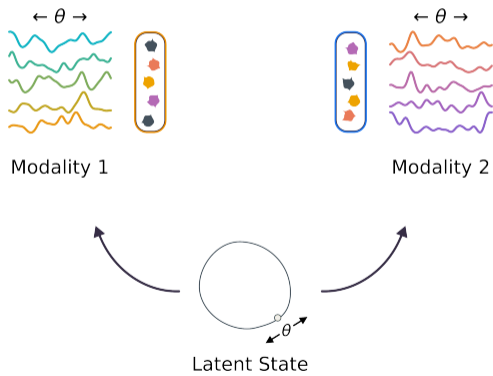


Modality 1

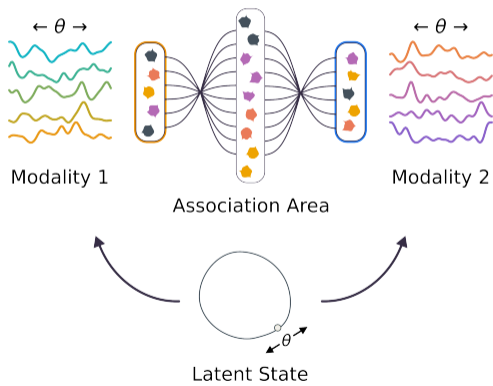
Modality 2



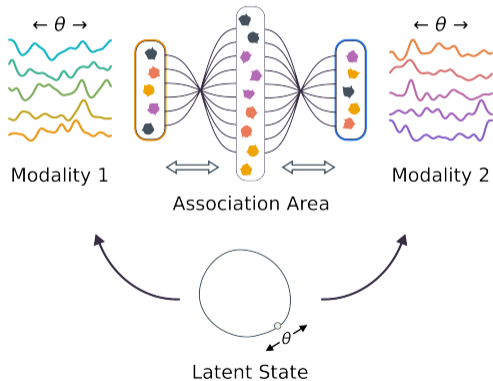
- Low-D latent variable
- Different areas, correlated variability



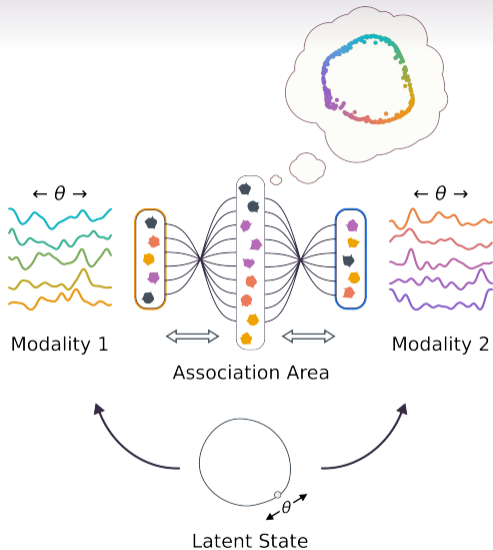
- Low-D latent variable
- Different areas, correlated variability



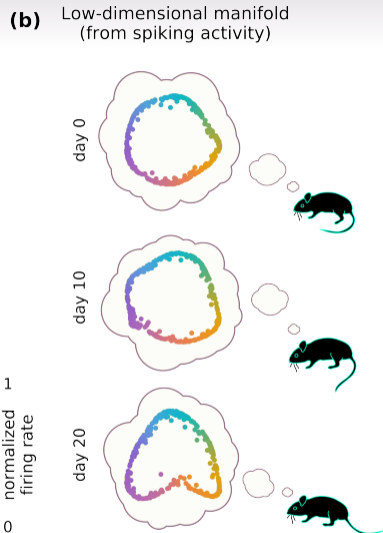
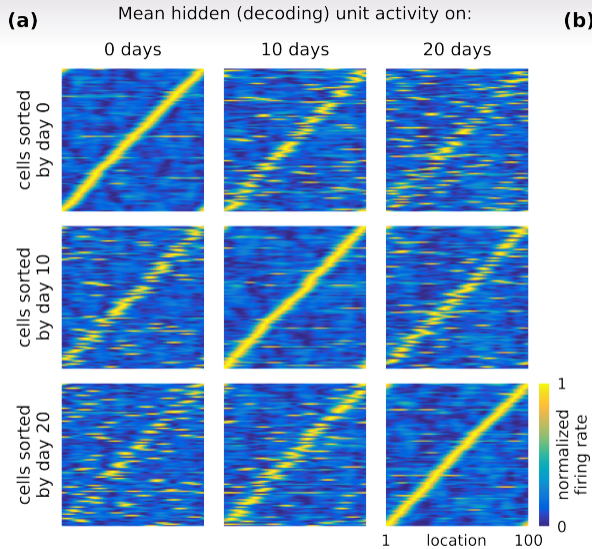
- Low-D latent variable
- Different areas, correlated variability



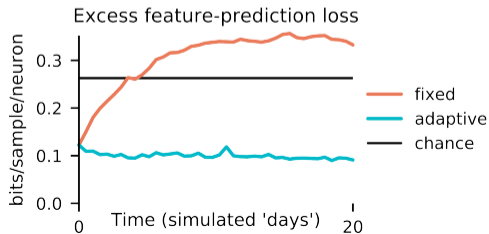
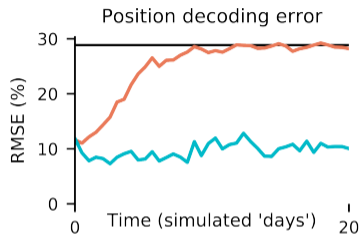
- Low-D latent variable
- Different areas, correlated variability
- Restricted Boltzmann Machine
 - Stochastic, binary
 - Non-negative activity



- Low-D latent variable
- Different areas, correlated variability
- Restricted Boltzmann Machine
 - Stochastic, binary
 - Non-negative activity
- Drift: noise to synaptic weights
 - Train continuously
 - Maintain mean rates
 - Normalize population responses (units compete)



Ongoing learning addresses drift



Looks like drift

- Code in "association" area changes
- Population low-D task structure stable

Consistent readout via unsupervised learning

- Correlated activation \rightarrow error feedback
- Plasticity: readouts to learn as quickly as the representation changes

Problem: Forgets easily

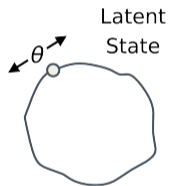
How could stable internal codes coexist with such unstable representations?

Q: How to achieve stable interpretations of unstable codes?

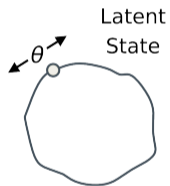
A: Homeostatic mechanisms create stability without error feedback.

Model Drift

Model encoding drift

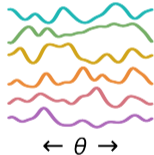


Model encoding drift

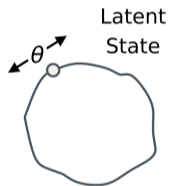


Input Features

$\mathbf{s}(\theta)$

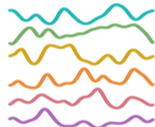


Model encoding drift



Input Features

$$\mathbf{s}(\theta)$$



$$\leftarrow \theta \rightarrow$$

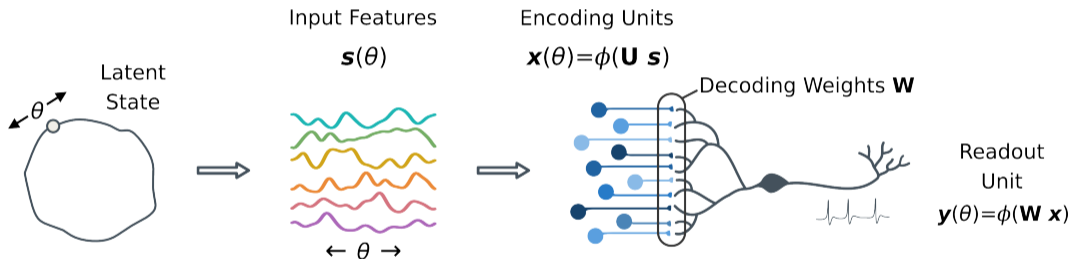


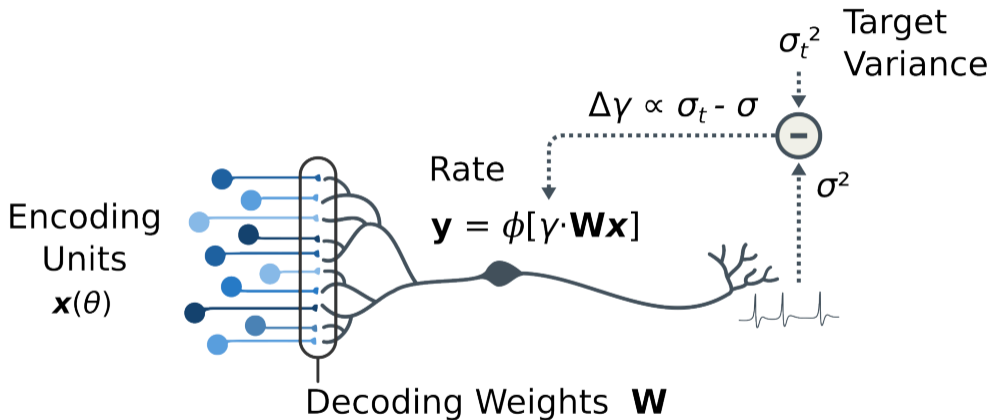
Encoding Units

$$\mathbf{x}(\theta) = \phi(\mathbf{U} \mathbf{s})$$



Model encoding drift

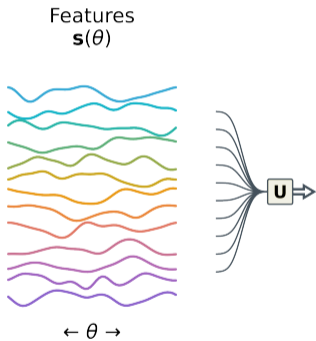




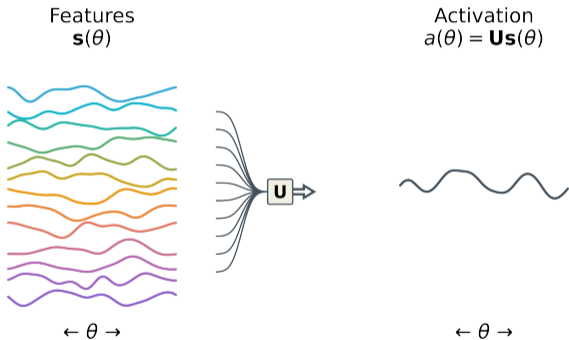
Neurons extract maxima of activation function over θ



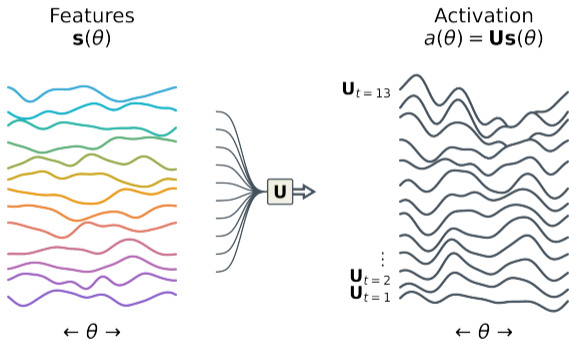
Neurons extract maxima of activation function over θ



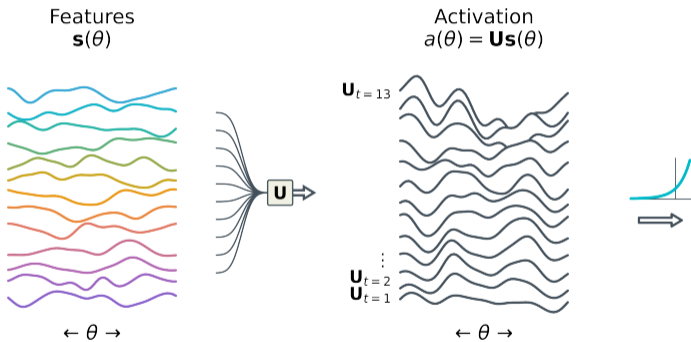
Neurons extract maxima of activation function over θ



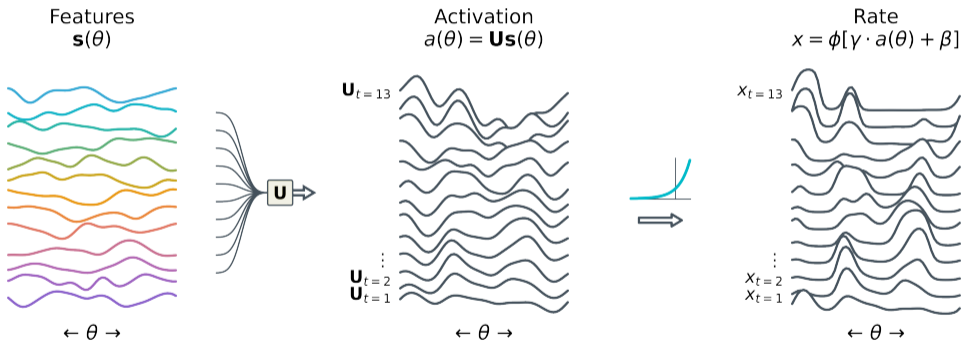
Neurons extract maxima of activation function over θ



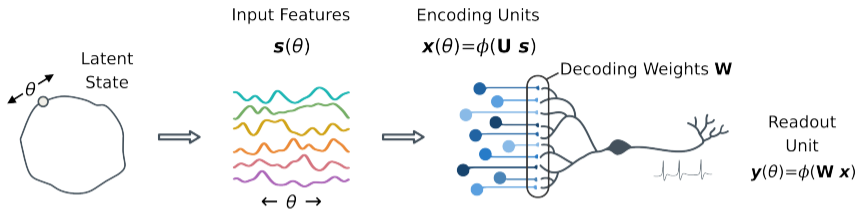
Neurons extract maxima of activation function over θ



Neurons extract maxima of activation function over θ



Hard to change preferred tuning



$$x(\theta) = \mathbf{U}\mathbf{s}(\theta), \quad \theta_0 = \underset{\theta}{\operatorname{argmax}}[x(\theta)]$$

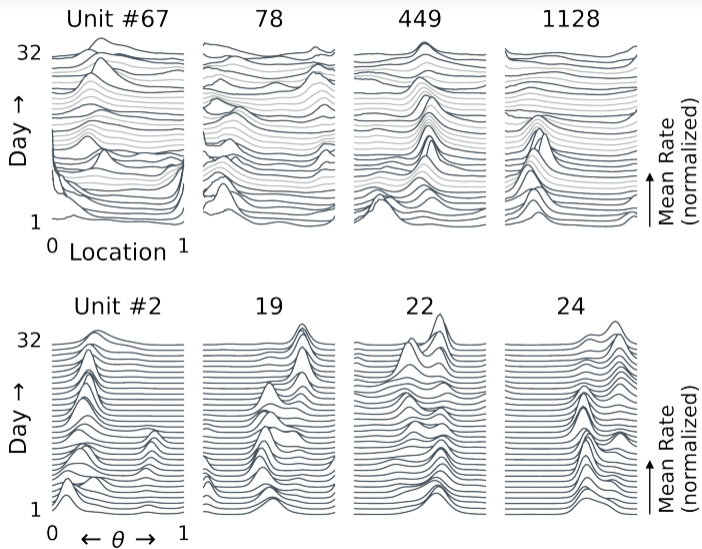
$$\Delta\theta_0 = - [\nabla_{\theta}\nabla_{\theta}^{\top}x(\theta)]^{-1} [\mathbf{U}\nabla_{\theta}\Delta\mathbf{s}(\theta)]$$

$\Delta\mathbf{s}(\theta)$ must resemble $\Delta\theta$ near the peak θ_0 , only $\nabla_{\theta}\Delta\mathbf{s}(\theta)$ matters

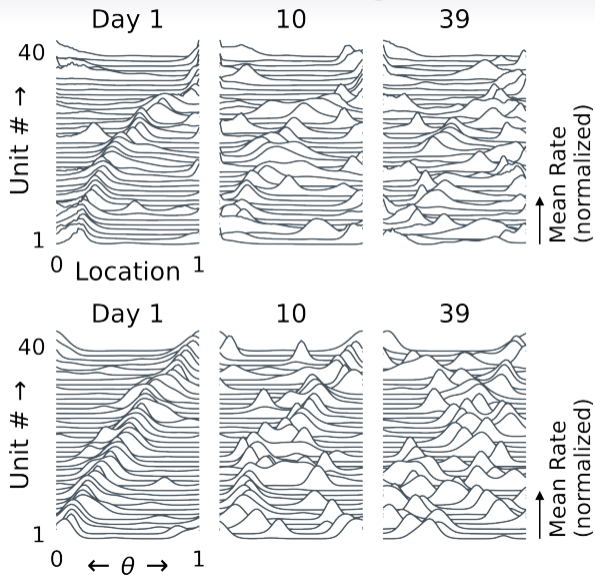
Any $\Delta\mathbf{s}(\theta)$ in null space of \mathbf{u} irrelevant

$[\nabla_{\theta}\nabla_{\theta}^{\top}x(\theta)]^{-1}$: sharper peaks are harder to move

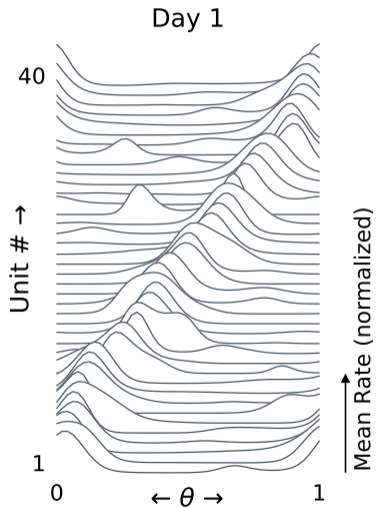
Model encoding drift



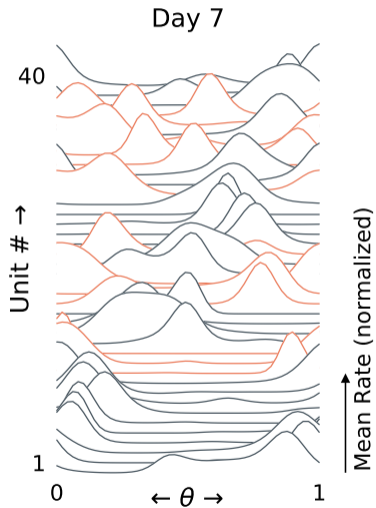
Model encoding drift



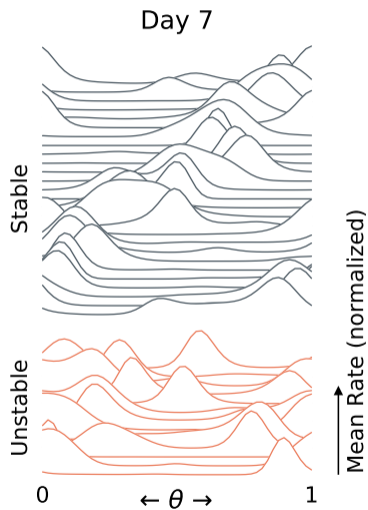
Drift is gradual



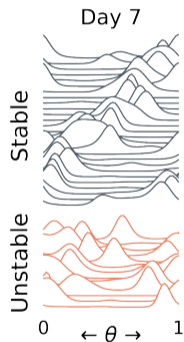
Drift is gradual



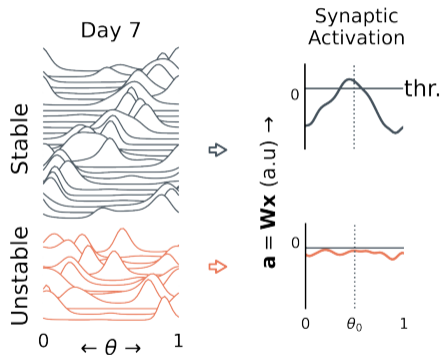
Drift is gradual



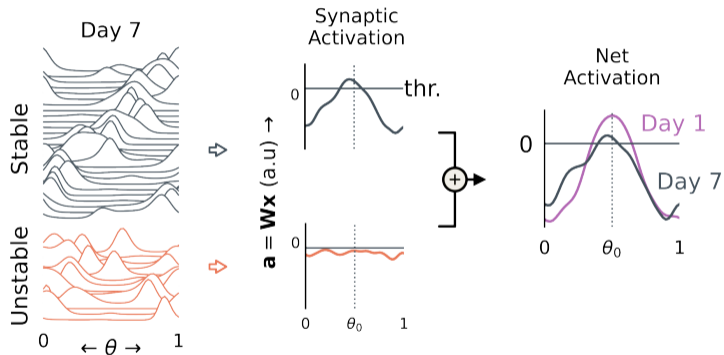
Drift causes loss of excitability, not tuning



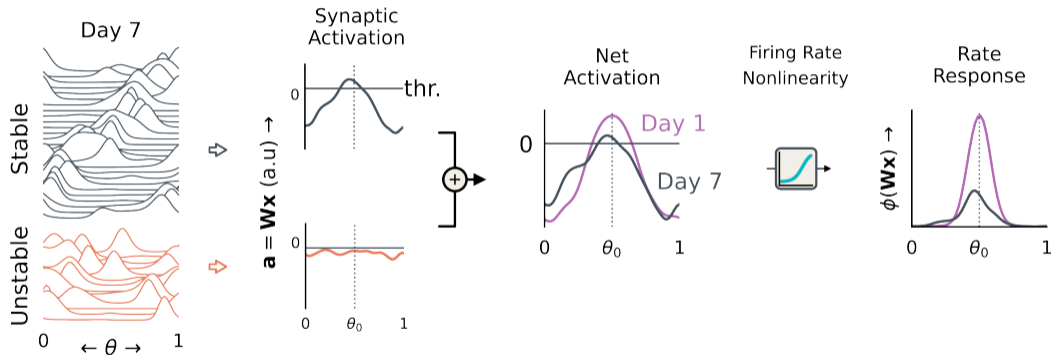
Drift causes loss of excitability, not tuning



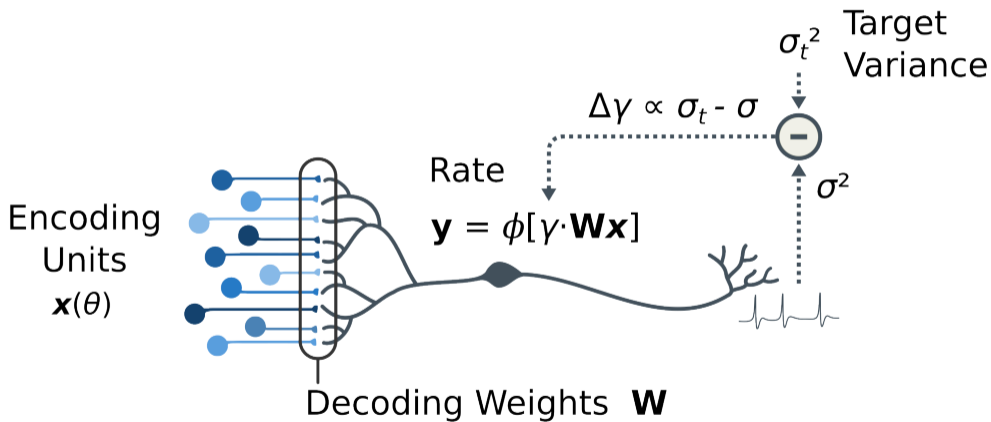
Drift causes loss of excitability, not tuning



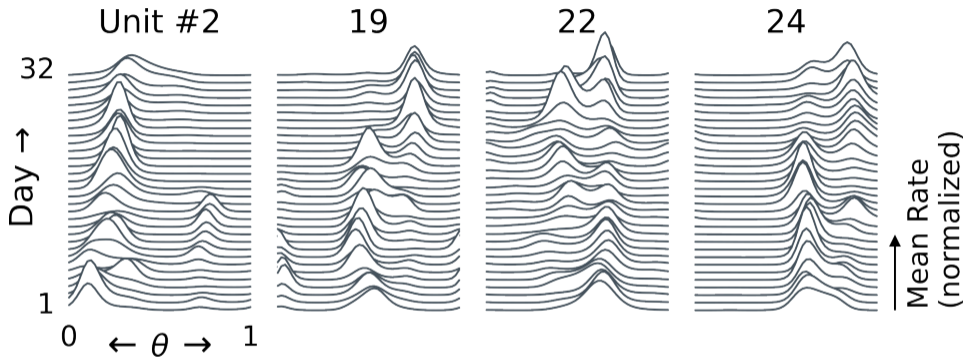
Drift causes loss of excitability, not tuning



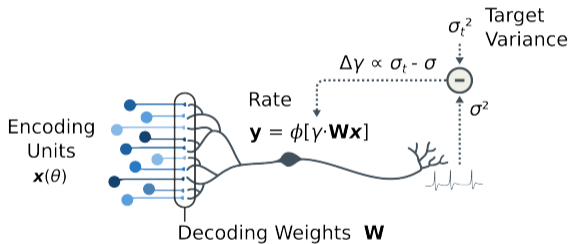
Sensitivity Homeostasis



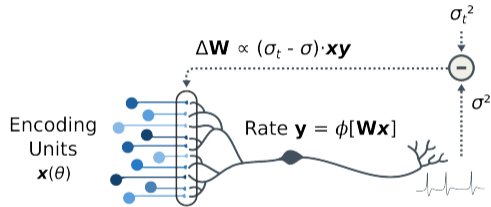
Sensitivity Homeostasis



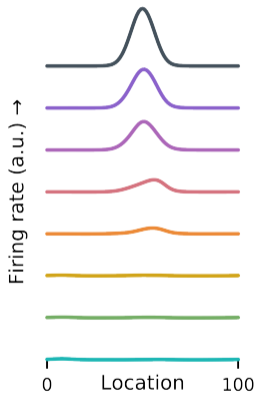
Sensitivity Homeostasis



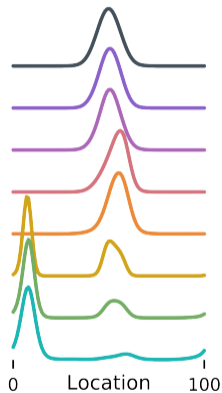
Hebbian Homeostasis



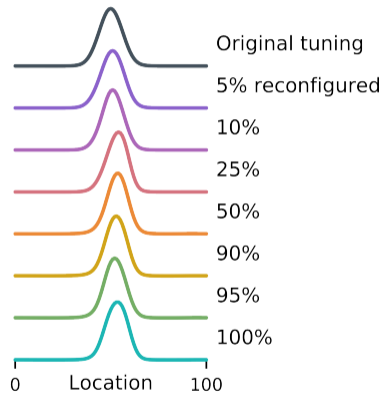
(a) Fixed Weights:
Loss of Excitability



(b) Sensitivity Homeostasis:
Punctuated Stability



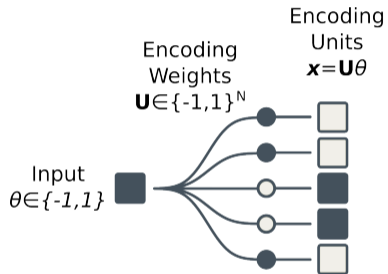
(c) Hebbian Homeostasis:
Stable Readout



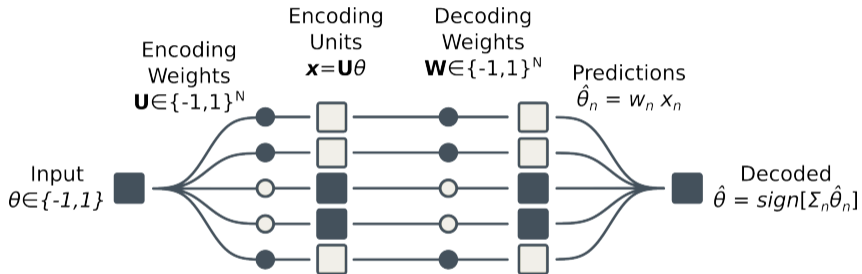
Why does this work?

Binary Threshold Analogy

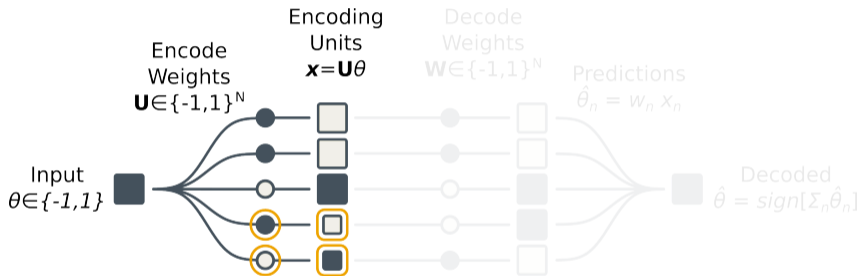
Error-Correcting Code



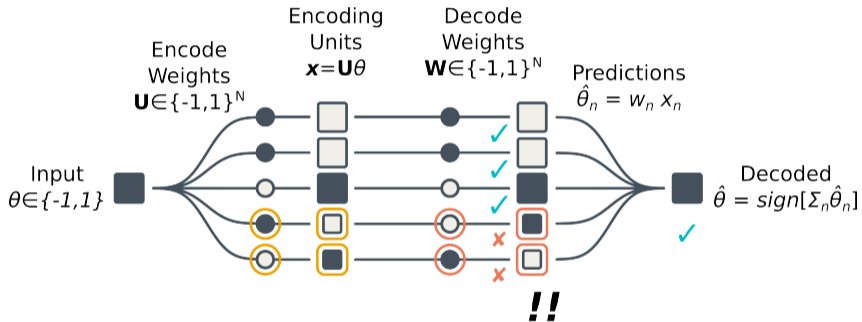
Error-Correcting Code



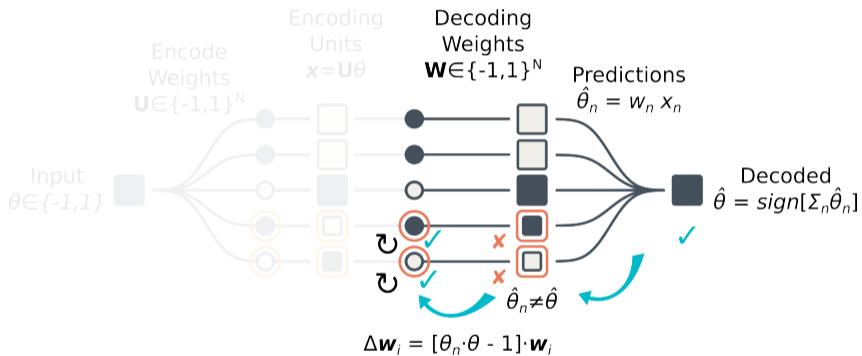
Error-Correcting Code



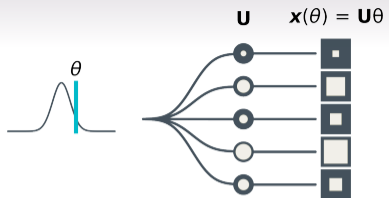
Error-Correcting Code



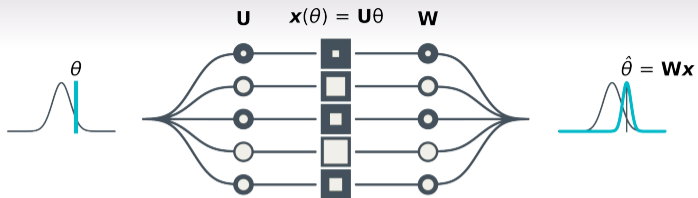
Self-Healing Code



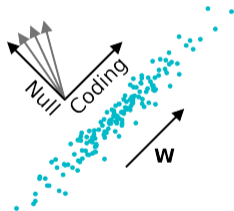
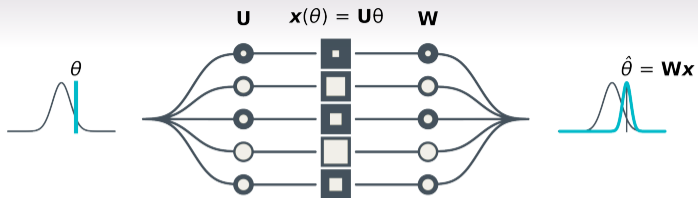
Linear Analogy



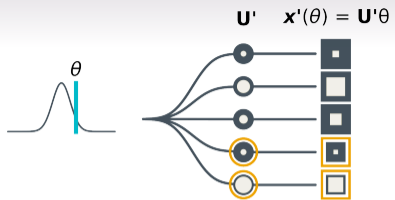
Redundant linear encoder



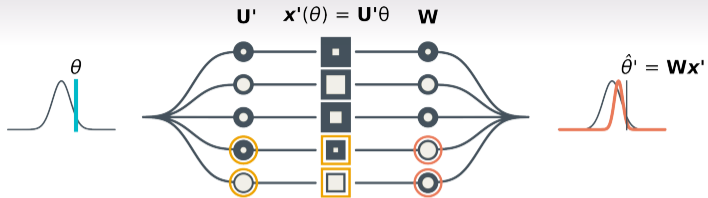
Low-D structure in High-D space



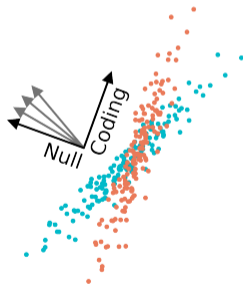
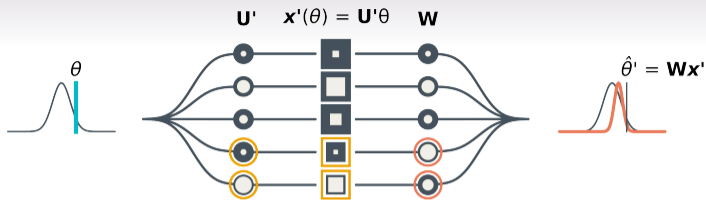
Many null-dimensions, weights align with signal dimensions



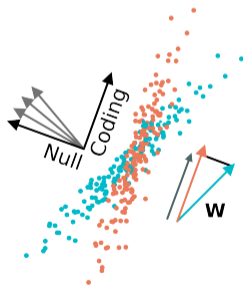
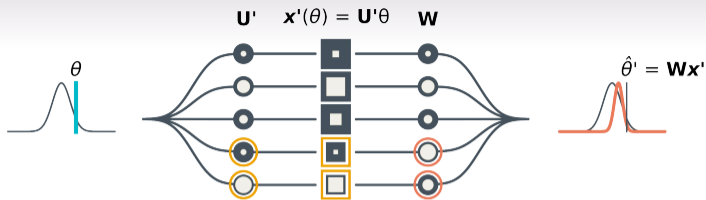
Small change in encoding...



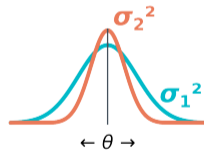
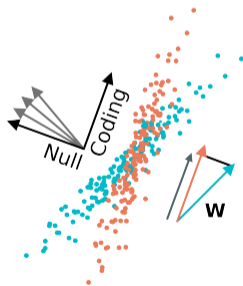
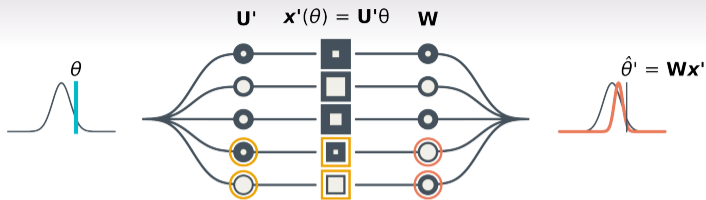
Loss of drive to readout



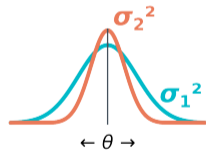
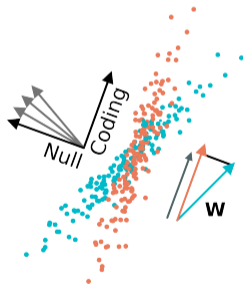
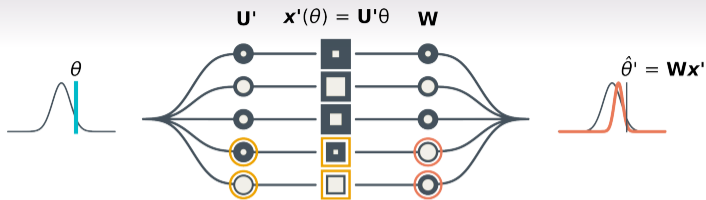
New embedding of low-D structure



Weights no longer match signal variability

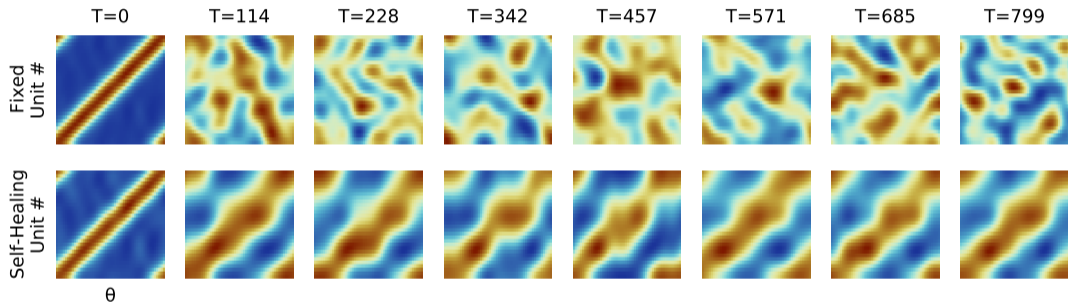


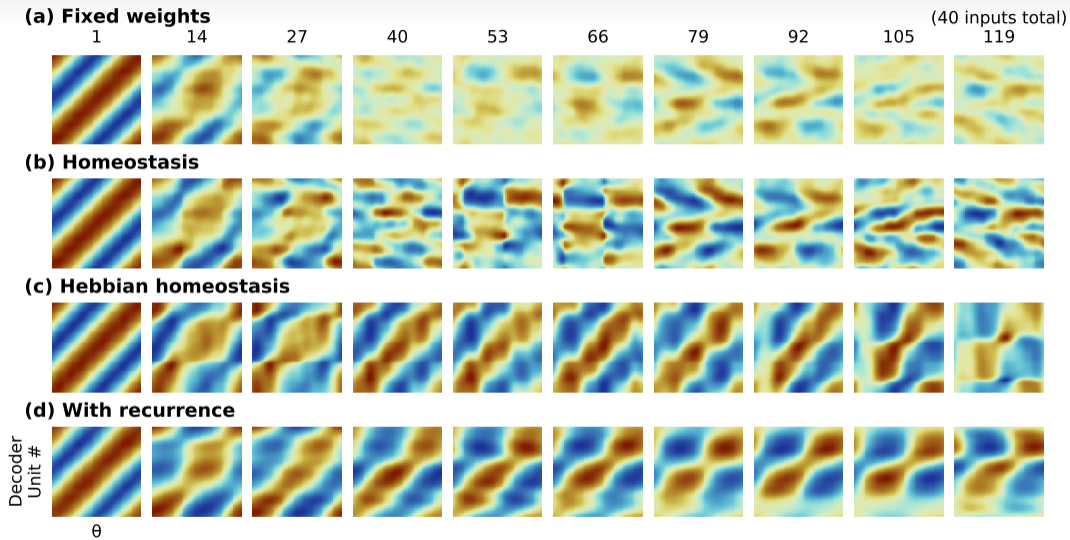
Detect loss of readout sensitivity



$$\Delta \mathbf{W}^T \propto (\sigma_1 - \sigma_2) \mathbf{x} \hat{\theta}' = (\sigma_1 - \sigma_2) \mathbf{x} \mathbf{x}^T \mathbf{W}^T$$

Hebbian homeostasis: realign weights to low-D structure





Linear-nonlinear readout

$$y(\theta) = \phi[\mathbf{W}\mathbf{x}(\theta)]$$

Sensitivity $y'(\theta) = \phi'[\mathbf{W}\mathbf{x}(\theta)]$

Drift $\Delta\mathbf{x}$; Average squared tuning change:

$$\langle \Delta_y^2 \rangle = \mathbf{W} \langle \Delta\mathbf{x}\Delta\mathbf{x}^T \cdot y'(\theta)^2 \rangle \mathbf{W}^T$$

Average sensitivity: $\|y'\|^2 = \int_{d\theta} y'(\theta)^2$

Normalized sensitivity: $\rho(\theta) = y'(\theta)^2 / \|y'\|^2$

$$\langle \Delta_y^2 \rangle = \|y'\|^2 \cdot \mathbf{W} \langle \Delta\mathbf{x}\Delta\mathbf{x}^T \cdot \rho(\theta) \rangle \mathbf{W}^T = \|y'\|^2 \cdot \mathbf{W} \Sigma_{\Delta\mathbf{x}}^{\rho(\theta)} \mathbf{W}^T$$

~ **Binary: saturating responses make $\|y'\|^2$ small**

~ **Linear: Locally-re-weighted input drift $\Sigma_{\Delta\mathbf{x}}^{\rho(\theta)}$ is low rank**

(remember)

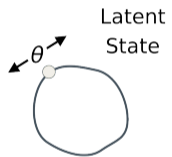
Representational drift is gradual (or null). It can be tracked via error feedback. Ongoing practice could provide this feedback, via prediction errors. However, this does not lead to stable internal representations.

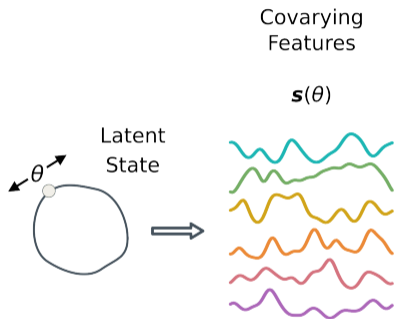
Model drift as shifting encoding weights

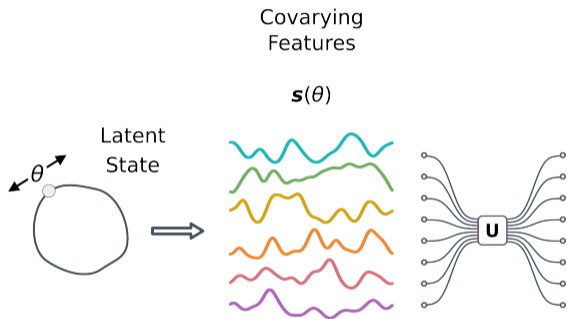
- Activity has low-D structure
- Sensitivity homeostasis leads to punctuated stability: occasional large shifts
- Hebbian homeostasis uses redundancy to re-learn weights as drift occurs
 - Binary: hard to change saturated responses
 - Linear: track low-D subspace
- Leads to stable readouts of unstable codes

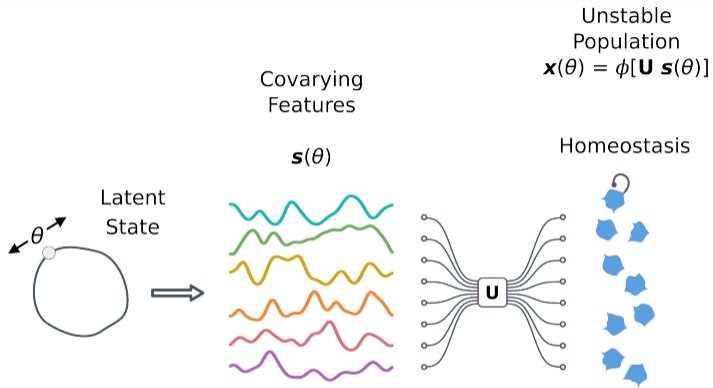
Next: Stabilizing population codes

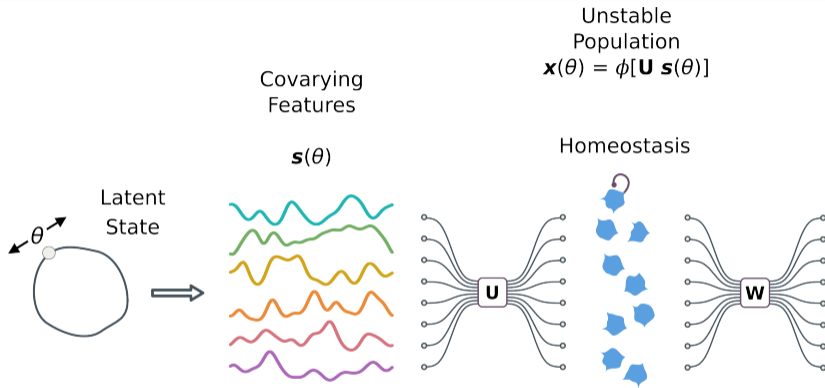
Population Interactions

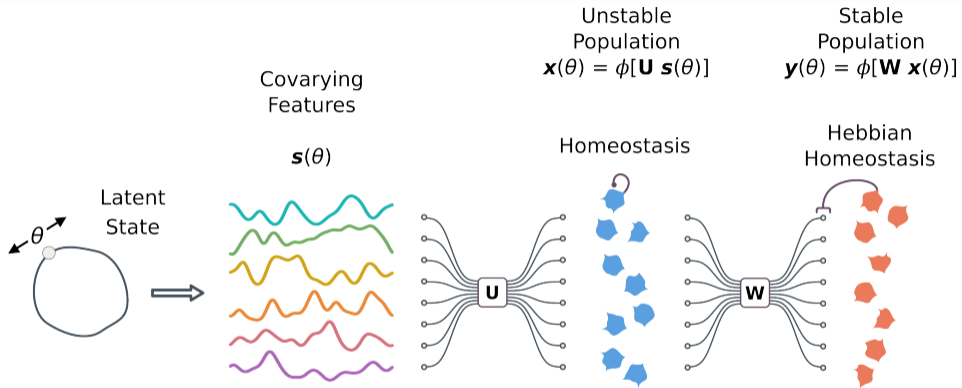


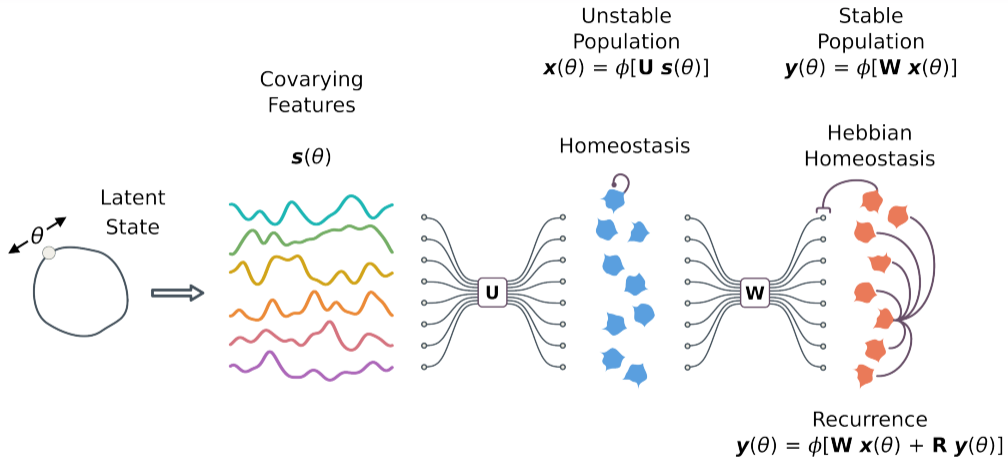


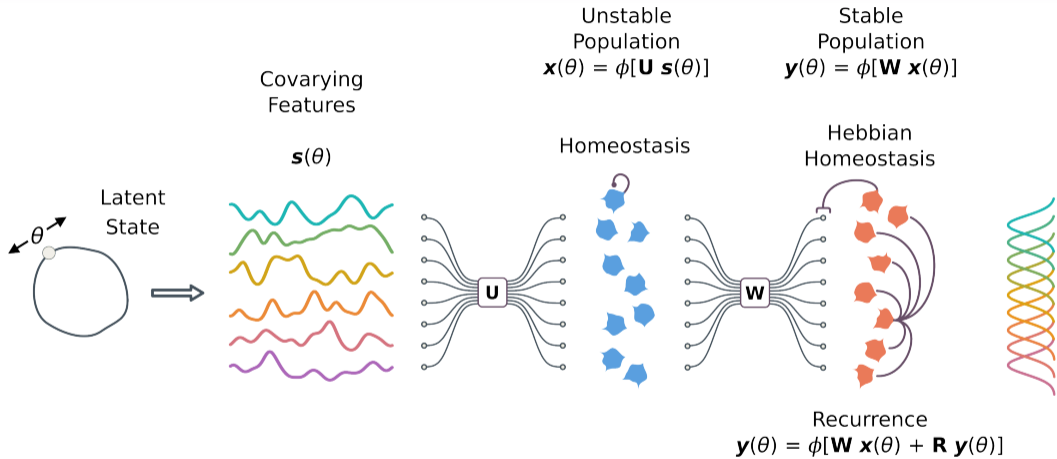












Population without Competition

Input Changes (60 inputs) →

0

24

48

72

96

120

Fixed Weights



Homeostasis



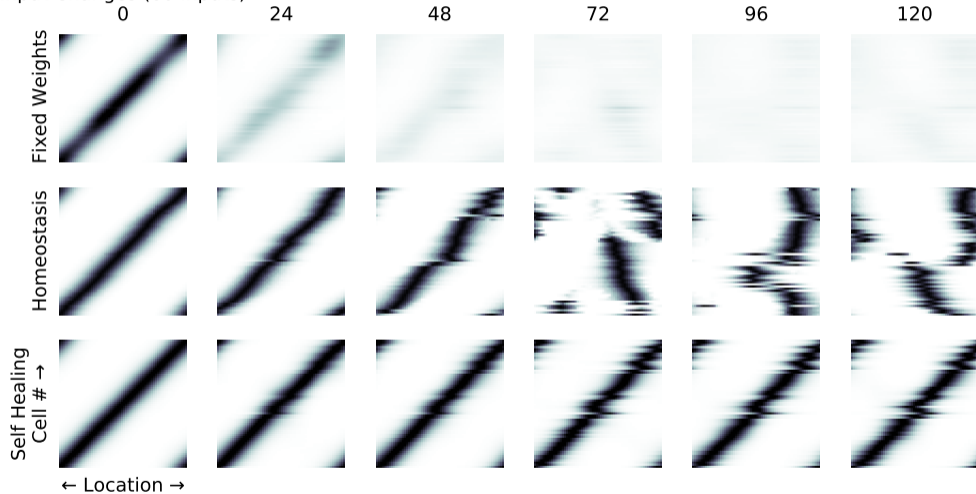
Self Healing
Cell # →



← Location →

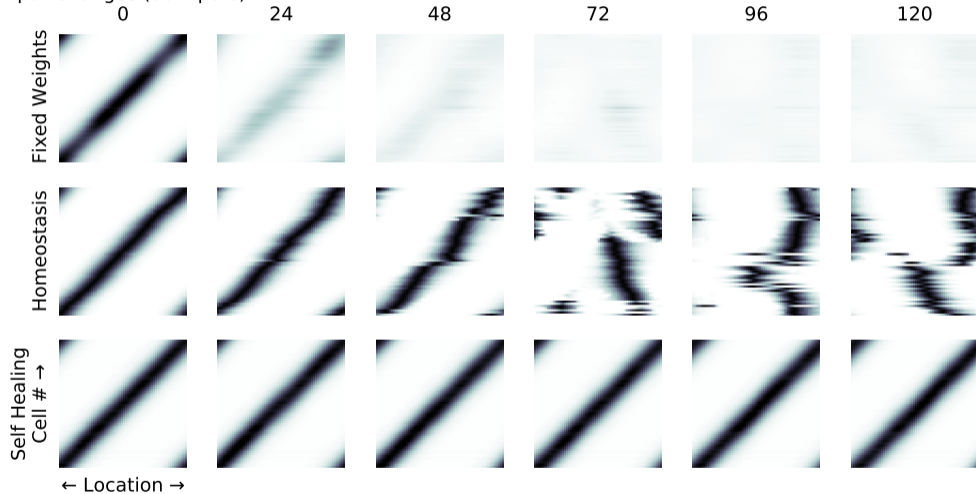
Population with Competition

Input Changes (60 inputs) →



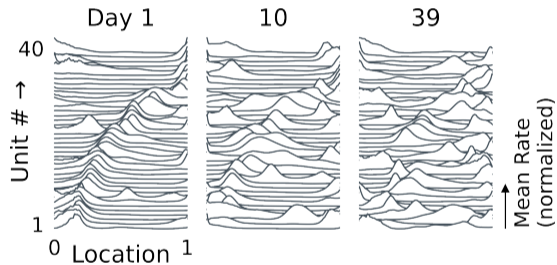
Population with Competition and Recurrence

Input Changes (60 inputs) →



(remember)

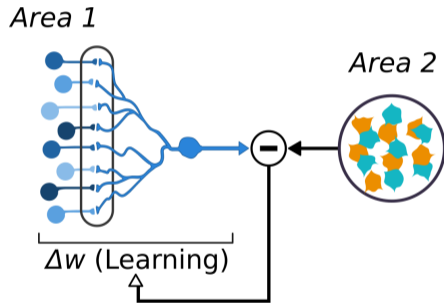
Drift is gradual/null:



(remember)

Drift is gradual/null:

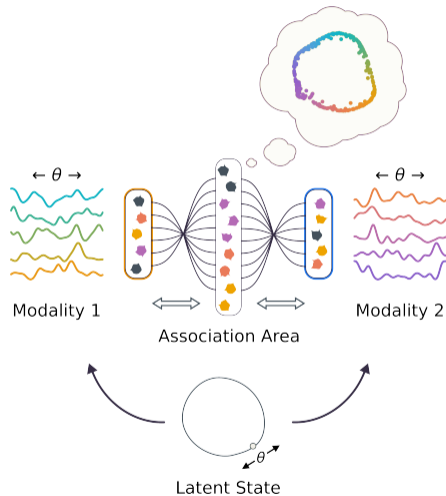
- Track with error feedback



(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

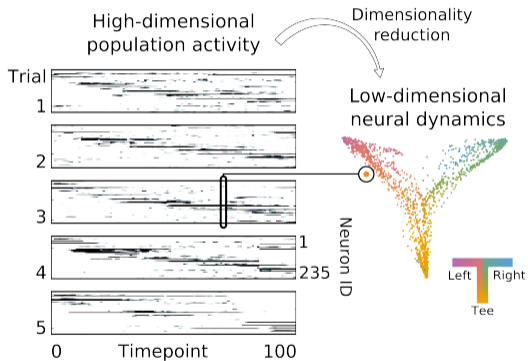


(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure



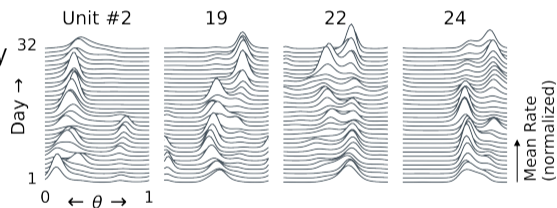
(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability



(remember)

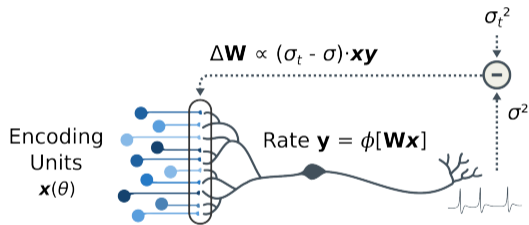
Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:



(remember)

Drift is gradual/null:

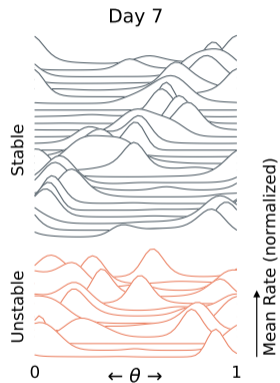
- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:

- Re-learn tuning as inputs change



(remember)

Drift is gradual/null:

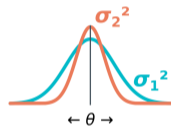
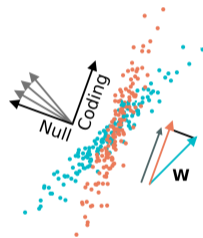
- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:

- Re-learn tuning as inputs change
- Binary: hard to change saturated responses
- Linear: track low-D subspace



$$\Delta \mathbf{W}^T \propto (\sigma_1 - \sigma_2) \mathbf{x} \theta^i = (\sigma_1 - \sigma_2) \mathbf{x} \mathbf{x}^T \mathbf{W}^T$$

(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:

- Re-learn tuning as inputs change
- Binary: hard to change saturated responses
- Linear: track low-D subspace

Population interactions

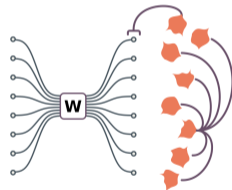
Unstable
Population
 $\mathbf{x}(\theta) = \phi[\mathbf{U} \mathbf{s}(\theta)]$

Stable
Population
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta)]$

Homeostasis



Hebbian
Homeostasis



Recurrence
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta) + \mathbf{R} \mathbf{y}(\theta)]$

(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:

- Re-learn tuning as inputs change
- Binary: hard to change saturated responses
- Linear: track low-D subspace

Population interactions

- Normalize: competition ensures coverage

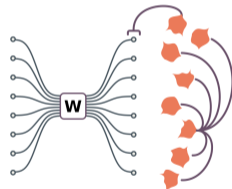
Unstable
Population
 $\mathbf{x}(\theta) = \phi[\mathbf{U} \mathbf{s}(\theta)]$

Stable
Population
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta)]$

Homeostasis



Hebbian
Homeostasis



Recurrence
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta) + \mathbf{R} \mathbf{y}(\theta)]$

(remember)

Drift is gradual/null:

- Track with error feedback
- Ongoing practice provides this

Model drift: Inputs have low-D structure

- Hard to change tuning; punctuated stability

Hebbian homeostasis:

- Re-learn tuning as inputs change
- Binary: hard to change saturated responses
- Linear: track low-D subspace

Population interactions

- Normalize: competition ensures coverage
- Recurrent connections → stable readout

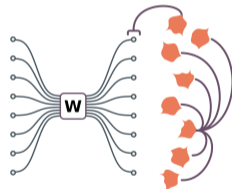
Unstable
Population
 $\mathbf{x}(\theta) = \phi[\mathbf{U} \mathbf{s}(\theta)]$

Stable
Population
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta)]$

Homeostasis



Hebbian
Homeostasis



Recurrence
 $\mathbf{y}(\theta) = \phi[\mathbf{W} \mathbf{x}(\theta) + \mathbf{R} \mathbf{y}(\theta)]$

End of Content

Thanks to:



Dhruva Raman



Timothy O'Leary



Chris Harvey



Laura Driscoll



Adrianna Loback



Fulvio Forni



Alon Rubin



Yaniv Ziv

Aspects of this work published in:

Rule ME, Loback AR, Raman DV, Driscoll L, Harvey CD, O'Leary T. 2020. Stable task information from an unstable neural population. *eLife*

Rule ME, O'Leary T, Harvey CD. 2019. Causes and consequences of representational drift. *Current opinion in neurobiology* 58:141–147

Funding:

This work was supported by the Human Frontier Science Program (RGY0069), ERC Starting Grant (StG FLEXNEURO 716643) and grants from the NIH (NS089521, MH107620, NS108410)