

Using Meta-Plasticity to Discover the Biophysics of Learning

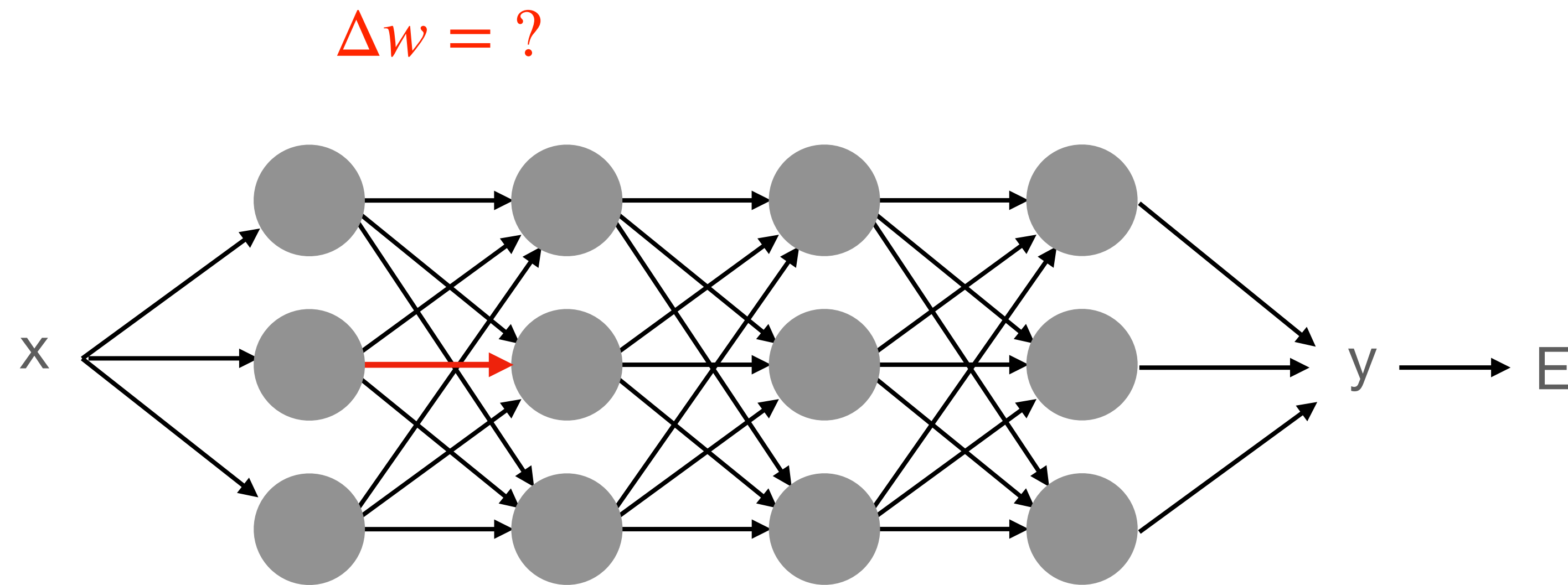
Navid Shervani-Tabar, Robert Rosenbaum

University of Notre Dame

Dept. of Applied and Computational Math and Stats



Learning in Artificial and Biological Neural Networks



Learning:

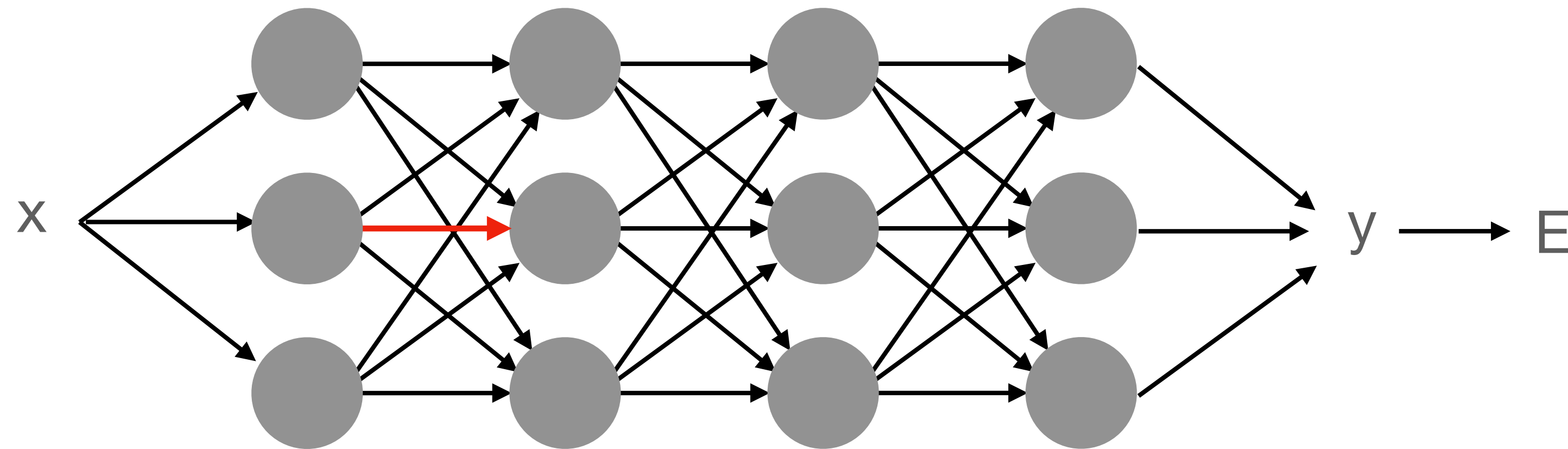
Update the weight in a way that decreases an error function.

Credit Assignment Problem:

How does a single weight contribute to the error?

Learning in Artificial and Biological Neural Networks

gradient descent: $\Delta w \propto -\frac{\partial E}{\partial w}$



Learning:

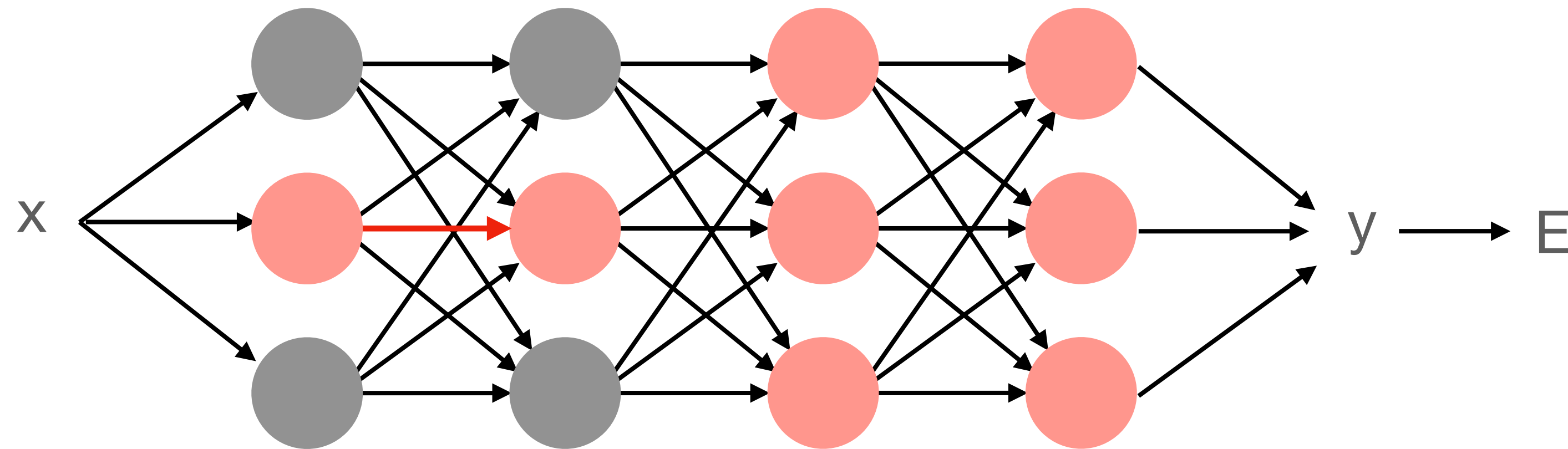
Update the weight in a way that decreases an error function.

Credit Assignment Problem:

How does a single weight contribute to the error?

Learning in Artificial and Biological Neural Networks

gradient descent: $\Delta w \propto -\frac{\partial E}{\partial w}$



Learning:

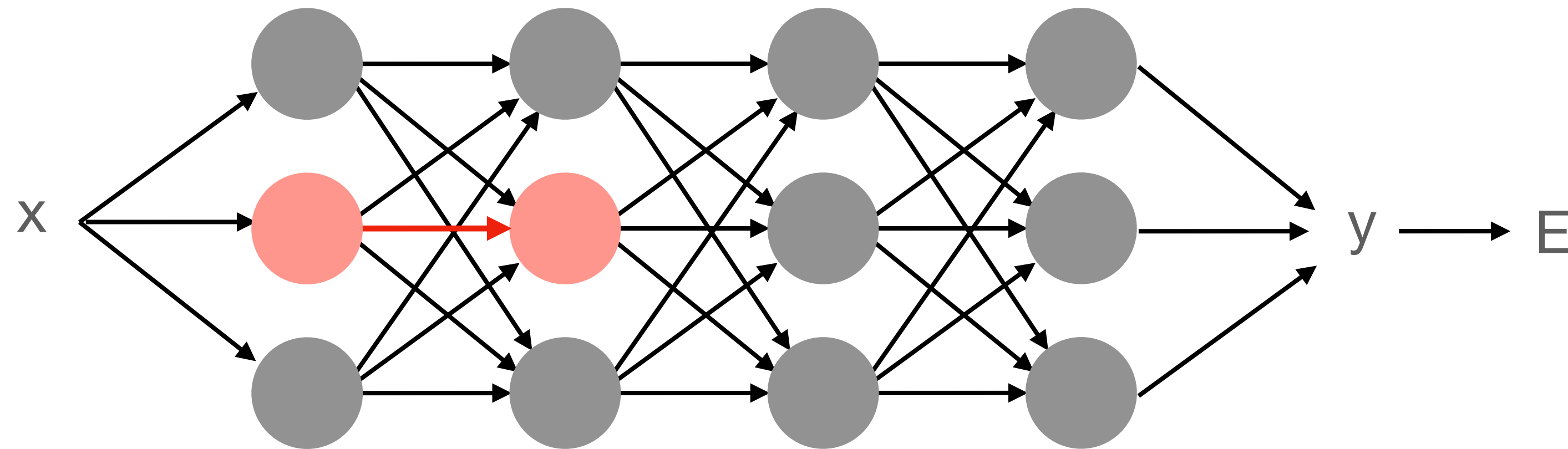
Update the weight in a way that decreases an error function.

Credit Assignment Problem:

How does a single weight contribute to the error?

Learning in Artificial and Biological Neural Networks

local plasticity: $\Delta w \approx F(v_{pre}, v_{post})$



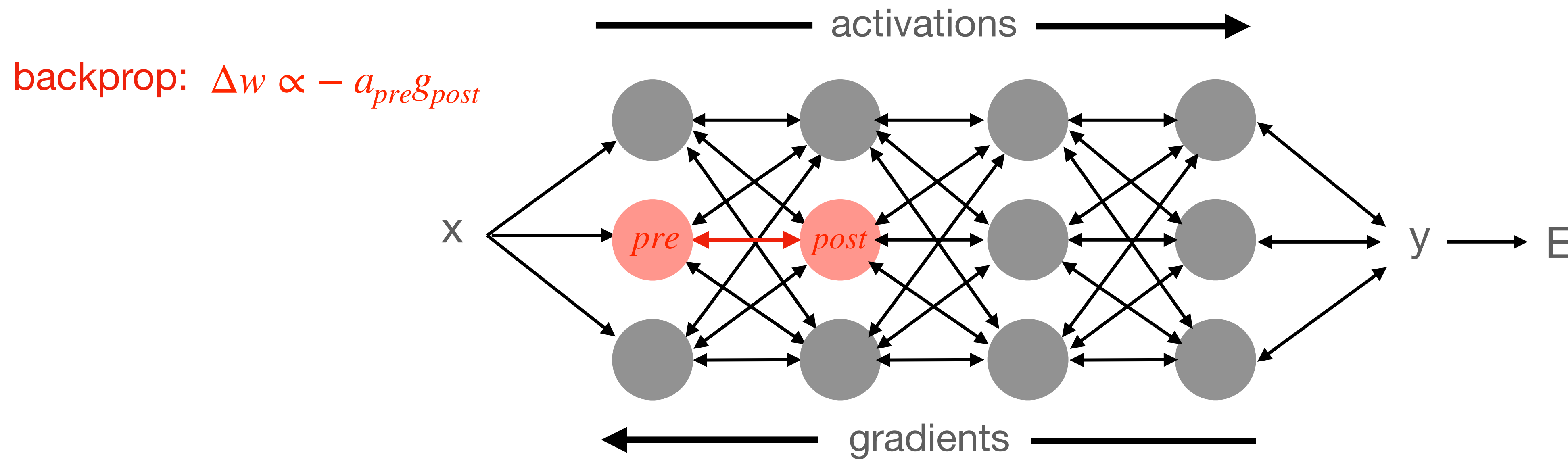
Local Plasticity

Changes to w are largely a function of pre- and post-synaptic activity.

Biological Credit Assignment Problem:

How to achieve effective learning under biologically relevant constraints?

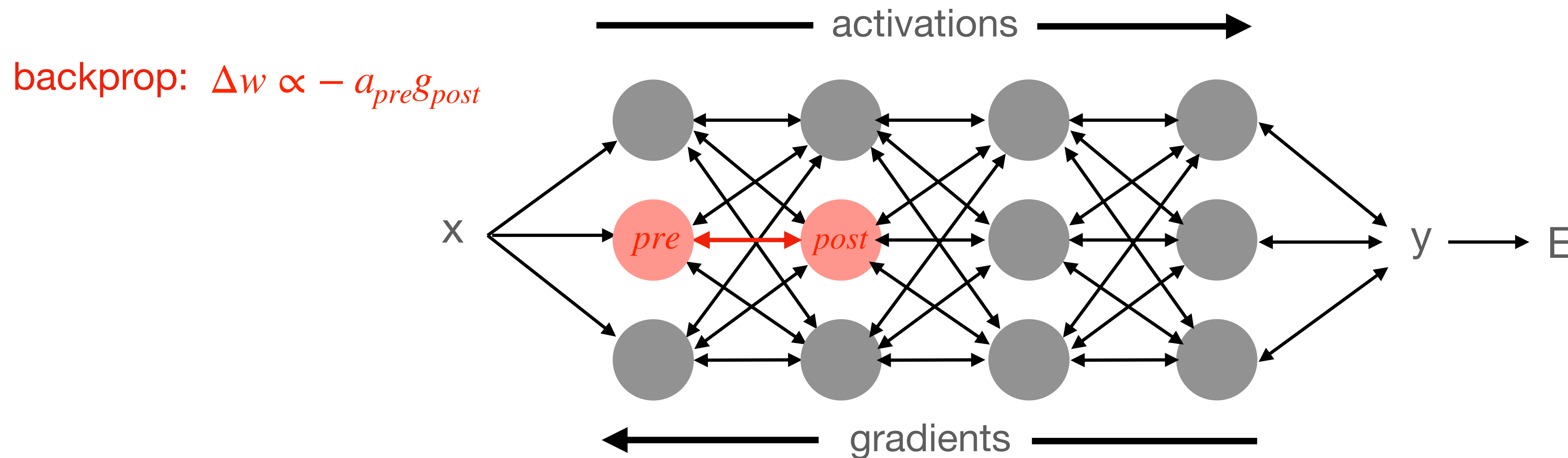
Learning in Artificial and Biological Neural Networks



Backpropagation:

Implements gradient descent with local learning rule by passing gradients backward.

Learning in Artificial and Biological Neural Networks



Article | [Published: 13 May 2021](#)

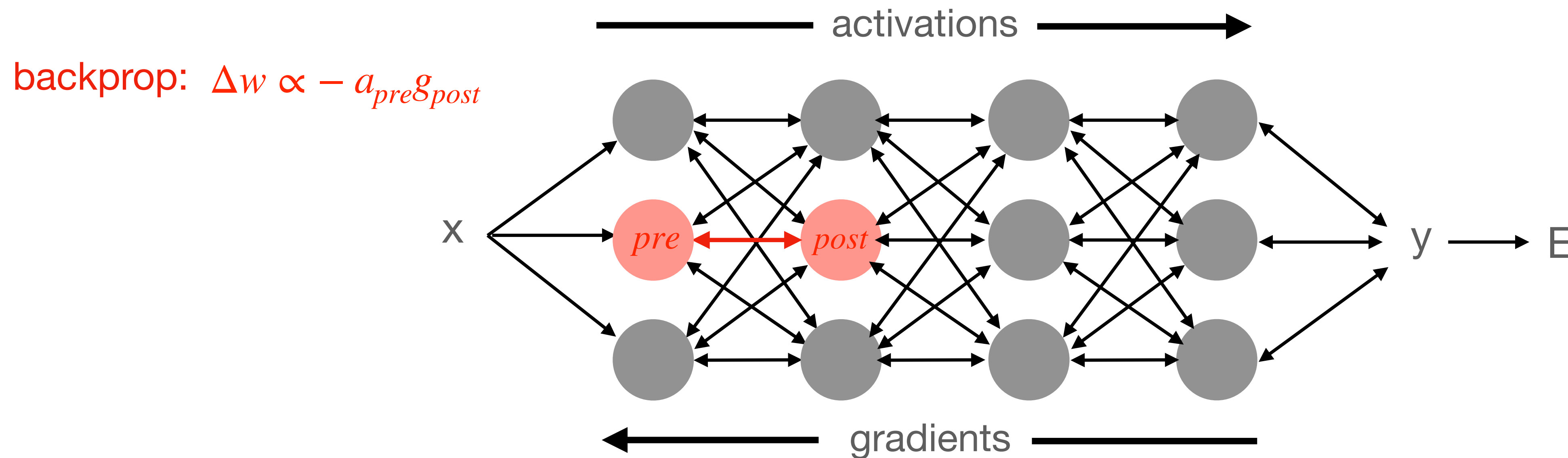
Burst-dependent synaptic plasticity can coordinate learning in hierarchical circuits

[Alexandre Payeur](#), [Jordan Guerguiev](#), [Friedemann Zenke](#), [Blake A. Richards](#)  & [Richard Naud](#) 

[Nature Neuroscience](#) **24**, 1010–1019 (2021) | [Cite this article](#)

activations and gradients
separately represented by
bursts and spikes?

Learning in Artificial and Biological Neural Networks



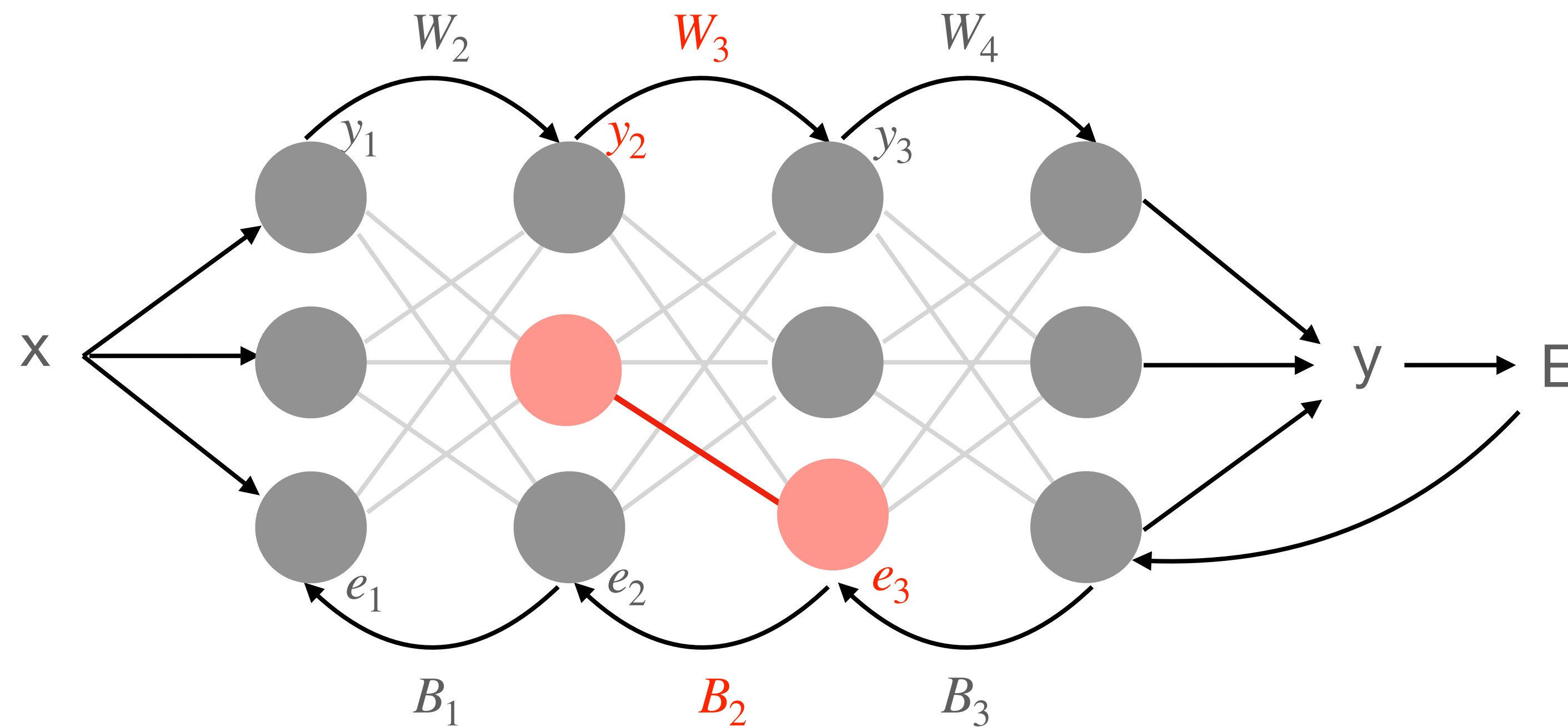
Backpropagation:

Implements gradient descent with local learning rule by passing gradients backward.

Problem:

Symmetric forward and backward connectivity.

Learning in Artificial and Biological Neural Networks



Backprop

$$\Delta W_{\ell}(j, k) = -\eta e_{\ell}(j) y_{\ell-1}(k)$$

$$B_{\ell} = W_{\ell+1}^T$$

“weight alignment”

Random feedback alignment

$$\Delta W_{\ell}(j, k) = -\eta e_{\ell}(j) y_{\ell-1}(k)$$

$$B_{\ell} = \text{fixed random}$$

(Lillicrap et al., 2016)

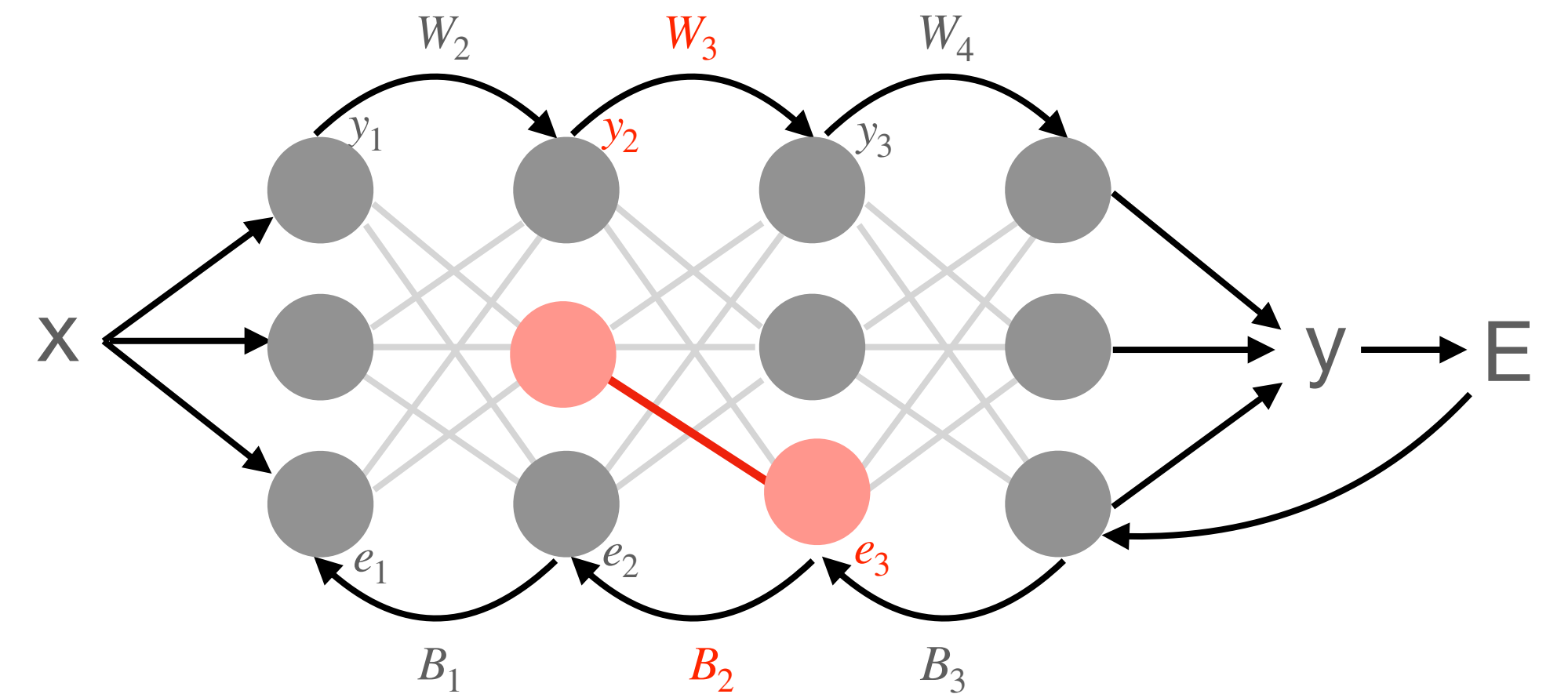
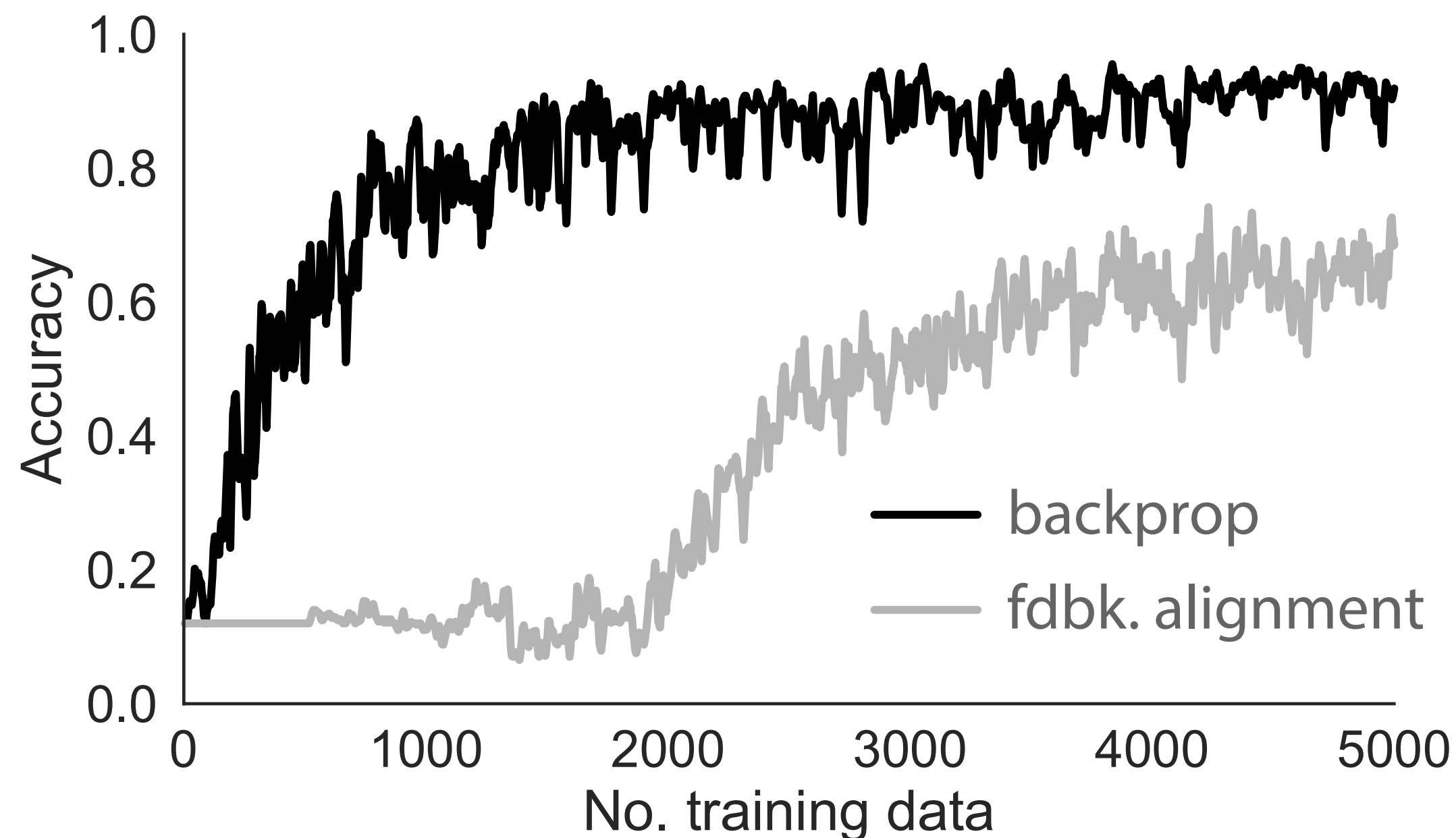
Learning in Artificial and Biological Neural Networks

Random feedback alignment

$$\Delta W_\ell(j, k) = -\eta e_\ell(j) y_{\ell-1}(k)$$

$B_\ell = \text{fixed random}$

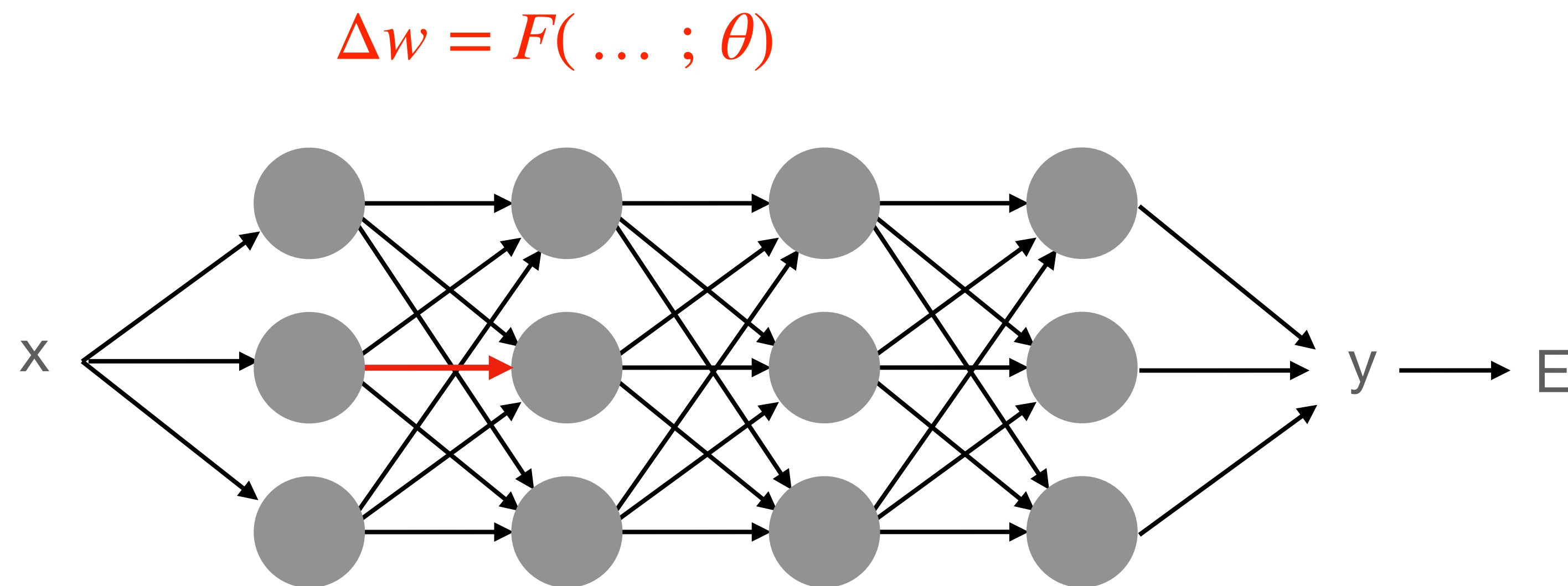
(Lillicrap et al., 2016)



Feedback alignment learns slowly and poorly in deep networks, especially with online learning (batch size = 1)

Meta-Plasticity to discover effective plasticity rules

- **Meta-Learning:** Optimize the learning rule across multiple tasks: “*learning to learn.*”
- **Meta-Plasticity:** Meta-learn a plasticity rule.



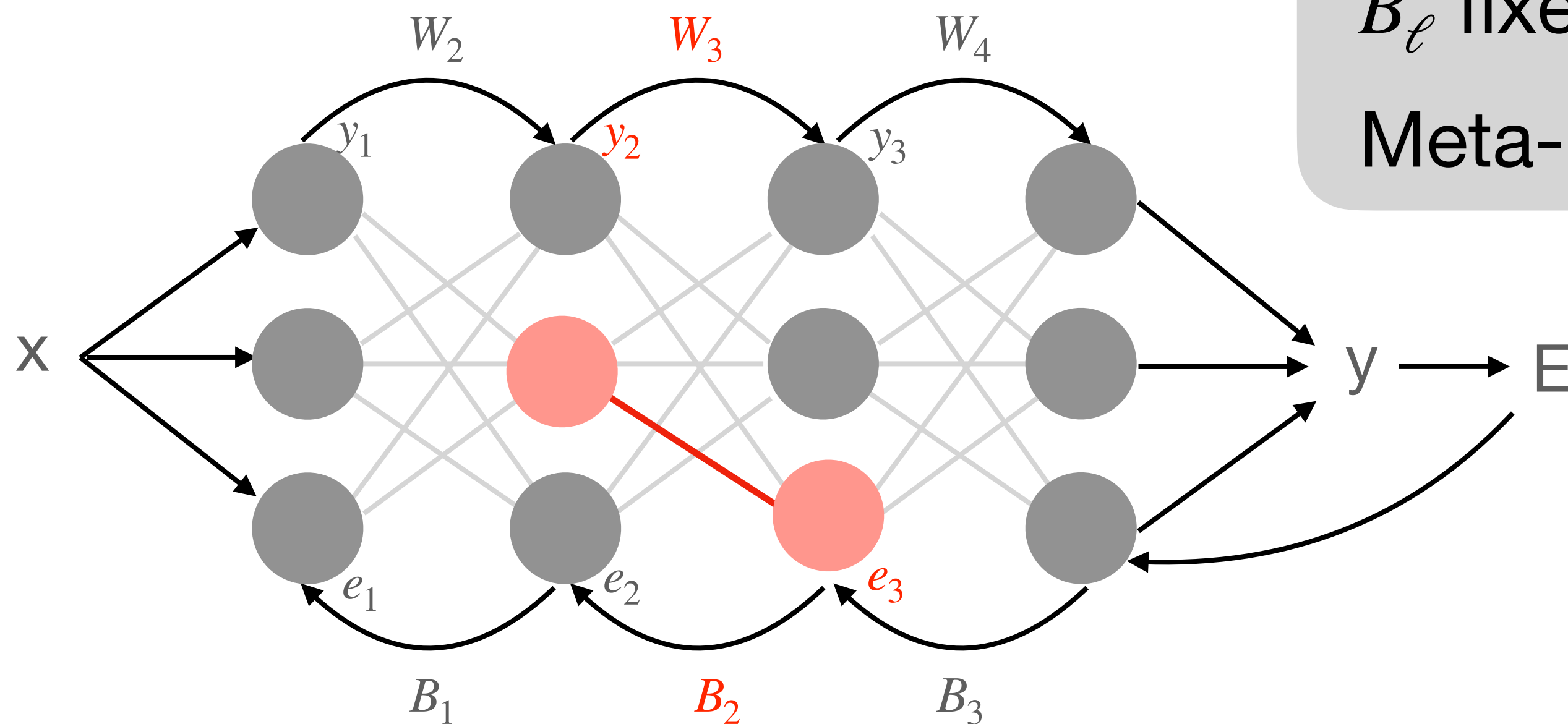
Meta-Plasticity to discover effective plasticity rules

- **Meta-Learning:** Optimize the learning rule across multiple tasks: “*learning to learn.*”
- **Meta-Plasticity:** Meta-learn a plasticity rule.
- Examples:
 - Meta-learn B_ℓ , initial W_ℓ , and parameters for plasticity rule at each synapse (Lindsey and Litwin-Kumar, 2020)
 - Meta-learn parameters for plasticity rule at each synapse. Retain the same W_ℓ across tasks (Miconi et al., 2019)
- **Results:** Better than backprop when generalizing to new tasks.
- **Problem:** Resulting learning rules are *not interpretable*, difficult to draw biological conclusions.

Our approach to meta-plasticity

- **Meta-parameter sharing**

- All synapses share the same local plasticity rule.
- Produces a single, *interpretable* plasticity rule.



$$\Delta W_{\ell}(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_{ℓ} fixed, random

Meta-learn θ

(See also: Confavreux et al., NeurIPS, 2020)

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

outer loop over tasks (“episodes”)

for each task:

Initialize W 's and B 's

for each data point:

Forward and backward pass to compute y 's and e 's

Update W 's using $\Delta W_\ell = F(\dots; \theta)$

Compute outputs \hat{Y} on unseen “query” data

Update theta using gradient of meta-cost: $J^{meta} = L(\hat{Y}, Y) + \lambda \|\theta\|_1$

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

for each task:

Initialize W 's and B 's

inner loop

for each data point:

Forward and backward pass to compute y 's and e 's

Update W 's using $\Delta W_\ell = F(\dots; \theta)$

Compute outputs \hat{Y} on unseen “query” data

Update theta using gradient of meta-cost: $J^{meta} = L(\hat{Y}, Y) + \lambda \|\theta\|_1$

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

for each task:

Initialize W 's and B 's inner loop learns from scratch

for each data point:

Forward and backward pass to compute y 's and e 's

Update W 's using $\Delta W_\ell = F(\dots; \theta)$

Compute outputs \hat{Y} on unseen “query” data

Update theta using gradient of meta-cost: $J^{meta} = L(\hat{Y}, Y) + \lambda \|\theta\|_1$

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

for each task:

Initialize W 's and B 's

for each data point: online learning (batch size 1)

Forward and backward pass to compute y 's and e 's

Update W 's using $\Delta W_\ell = F(\dots; \theta)$

Compute outputs \hat{Y} on unseen “query” data

Update theta using gradient of meta-cost: $J^{meta} = L(\hat{Y}, Y) + \lambda \|\theta\|_1$

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

for each task:

Initialize W 's and B 's

for each data point:

Forward and backward pass to compute y 's and e 's

Update W 's using $\Delta W_\ell = F(\dots; \theta)$

Compute outputs \hat{Y} on unseen “query” data

L1 meta-regularization
promotes simple plasticity rules

Update theta using gradient of meta-cost: $J^{meta} = L(\hat{Y}, Y) + \lambda \|\theta\|_1$

Our approach to meta-plasticity

$$\Delta W_\ell(j, k) = F(y_{pre}(k), y_{post}(j), e_{pre}(k), e_{post}(j); \theta)$$

B_ℓ fixed, random

Meta-learn θ

- 5-layer fully connected perceptron
- Each task is a subset of EMNIST
- Inner loop data set size = 256
- Linear combo of 10 bio-inspired plasticity rules:

$$F(\dots; \theta) = \sum_k \theta_k \mathcal{F}^{(k)}(\dots)$$

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_\ell \mathbf{y}_{\ell-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_\ell \mathbf{e}_{\ell-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_\ell \mathbf{e}_{\ell-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{\ell-1, \ell},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_\ell \mathbf{e}_{\ell-1}^T,$$

$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_\ell \mathbf{1}_\ell^T \mathbf{y}_\ell \mathbf{y}_{\ell-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_\ell \mathbf{y}_\ell^T \mathbf{W}_{\ell-1, \ell} \mathbf{e}_{\ell-1} \mathbf{e}_{\ell-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_\ell \mathbf{y}_\ell^T \mathbf{W}_{\ell-1, \ell} \mathbf{e}_{\ell-1} \mathbf{y}_{\ell-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_\ell \mathbf{y}_{\ell-1}^T \mathbf{W}_{\ell-1, \ell}^T \mathbf{e}_\ell \mathbf{e}_{\ell-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_\ell \mathbf{y}_{\ell-1}^T - (\mathbf{y}_\ell \mathbf{y}_\ell^T) \mathbf{W}_{\ell-1, \ell})$$

Results

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{l-1,l},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l \mathbf{e}_{l-1}^T,$$

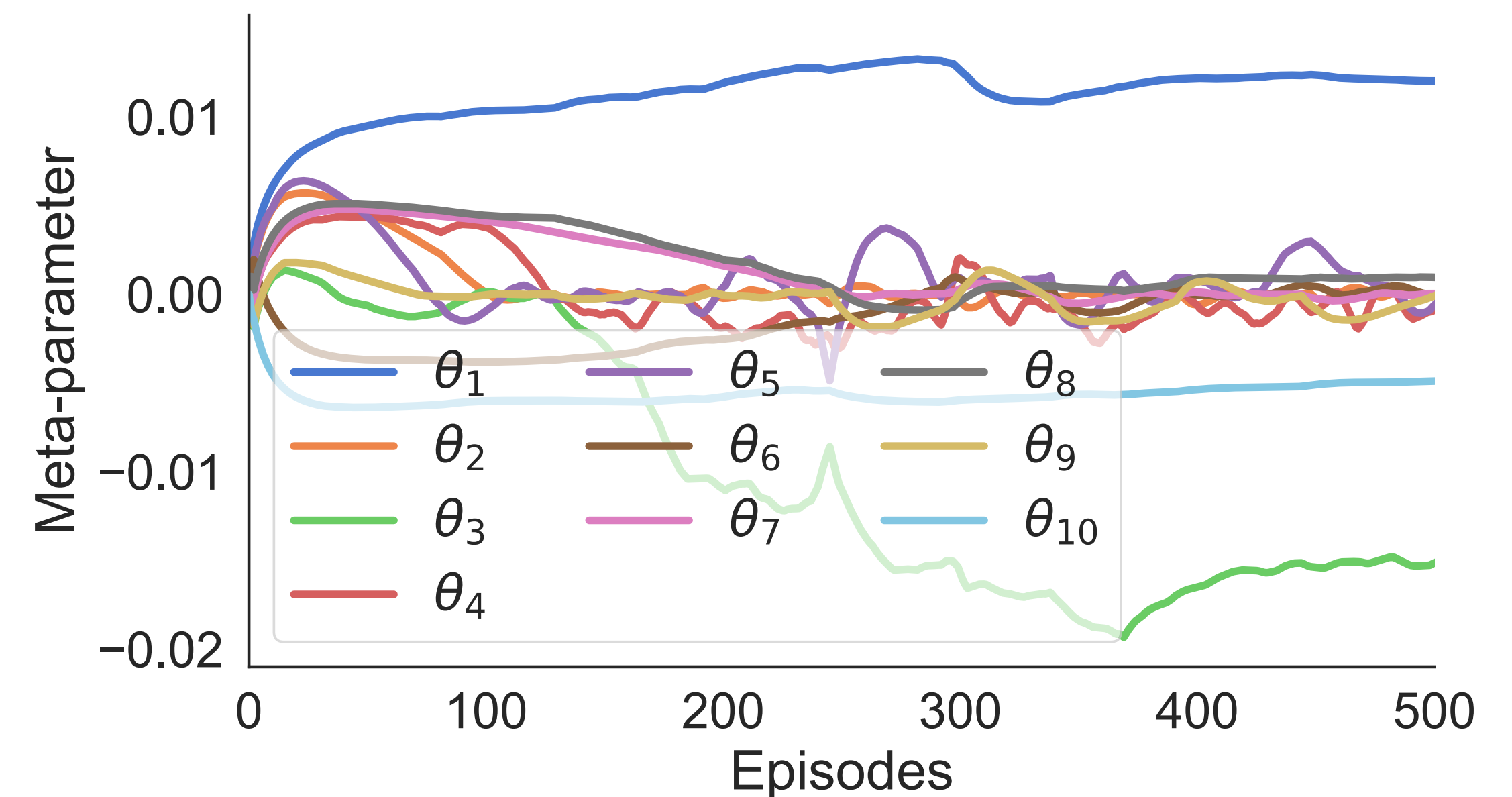
$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_l \mathbf{1}_l^T \mathbf{y}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_l \mathbf{y}_{l-1}^T \mathbf{W}_{l-1,l}^T \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_l \mathbf{y}_{l-1}^T - (\mathbf{y}_l \mathbf{y}_l^T) \mathbf{W}_{l-1,l})$$



Results

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{l-1,l},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_l \mathbf{1}_l^T \mathbf{y}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_l \mathbf{y}_{l-1}^T \mathbf{W}_{l-1,l}^T \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_l \mathbf{y}_{l-1}^T - (\mathbf{y}_l \mathbf{y}_l^T) \mathbf{W}_{l-1,l})$$

Oja's rule on activations.
Causes hidden layers to perform
PCA-like orthonormalization.

Results

$$\mathcal{F}^{(1)} = -\theta_1 e_l y_l^T$$

$$\mathcal{F}^{(2)} = -\theta_2 y_l e_l^T$$

$$\mathcal{F}^{(3)} = -\theta_3 e_l e_l^T$$

$$\mathcal{F}^{(4)} = -\theta_4 W_{l-1,l}$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l e_l^T$$

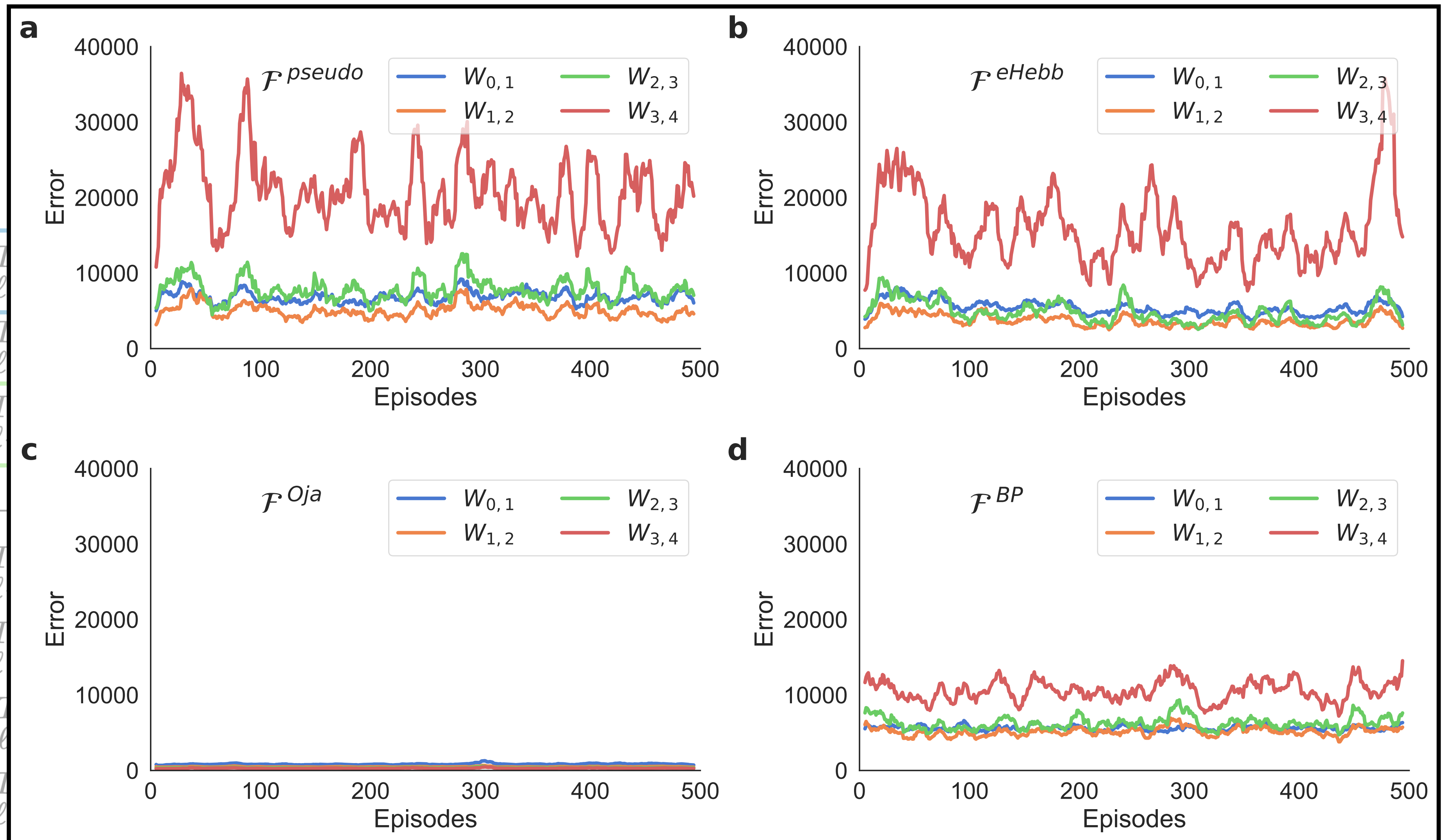
$$\mathcal{F}^{(6)} = -\theta_6 e_l \mathbf{1}_l^T$$

$$\mathcal{F}^{(7)} = -\theta_7 y_l y_l^T$$

$$\mathcal{F}^{(8)} = -\theta_8 e_l y_l^T$$

$$\mathcal{F}^{(9)} = -\theta_9 y_l y_{l-1}^T W_{l-1,l}^T e_l e_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (y_l y_{l-1}^T - (y_l y_l^T) W_{l-1,l})$$



Oja's rule on activations.
Causes hidden layers to perform
PCA-like orthonormalization.

Results

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{l-1,l},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_l \mathbf{1}_l^T \mathbf{y}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_l \mathbf{y}_{l-1}^T \mathbf{W}_{l-1,l}^T \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_l \mathbf{y}_{l-1}^T - (\mathbf{y}_l \mathbf{y}_l^T) \mathbf{W}_{l-1,l})$$

Pseudo-gradient term. If $\mathbf{W}_\ell = \mathbf{B}_{\ell-1}^T$, this term would implement backprop.

Oja's rule on activations.
Causes hidden layers to perform PCA-like orthonormalization.

Results

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{l-1,l},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_l \mathbf{1}_l^T \mathbf{y}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_l \mathbf{y}_{l-1}^T \mathbf{W}_{l-1,l}^T \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_l \mathbf{y}_{l-1}^T - (\mathbf{y}_l \mathbf{y}_l^T) \mathbf{W}_{l-1,l})$$

Pseudo-gradient term. If $W_\ell = B_{\ell-1}^T$, this term would implement backprop.

Hebbian plasticity on errors.
Theorem: In a linear network, this term pushes W_ℓ toward $B_{\ell-1}^T$

Oja's rule on activations.
Causes hidden layers to perform PCA-like orthonormalization.

Results

$$\mathcal{F}^{(1)} = -\theta_1 \mathbf{e}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(2)} = -\theta_2 \mathbf{y}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(3)} = -\theta_3 \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(4)} = -\theta_4 \mathbf{W}_{l-1,l},$$

$$\mathcal{F}^{(5)} = -\theta_5 \mathbf{1}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(6)} = -\theta_6 \mathbf{e}_l \mathbf{1}_l^T \mathbf{y}_l \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(7)} = -\theta_7 \mathbf{y}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(8)} = -\theta_8 \mathbf{e}_l \mathbf{y}_l^T \mathbf{W}_{l-1,l} \mathbf{e}_{l-1} \mathbf{y}_{l-1}^T,$$

$$\mathcal{F}^{(9)} = -\theta_9 \mathbf{y}_l \mathbf{y}_{l-1}^T \mathbf{W}_{l-1,l}^T \mathbf{e}_l \mathbf{e}_{l-1}^T,$$

$$\mathcal{F}^{(10)} = -\theta_{10} (\mathbf{y}_l \mathbf{y}_{l-1}^T - (\mathbf{y}_l \mathbf{y}_l^T) \mathbf{W}_{l-1,l})$$

Hebbian plasticity on errors.

Theorem: In a linear network, this term pushes W_ℓ toward $B_{\ell-1}^T$

Could this rule help align weights for learning via burst-dependent plasticity?

Article | [Published: 13 May 2021](#)

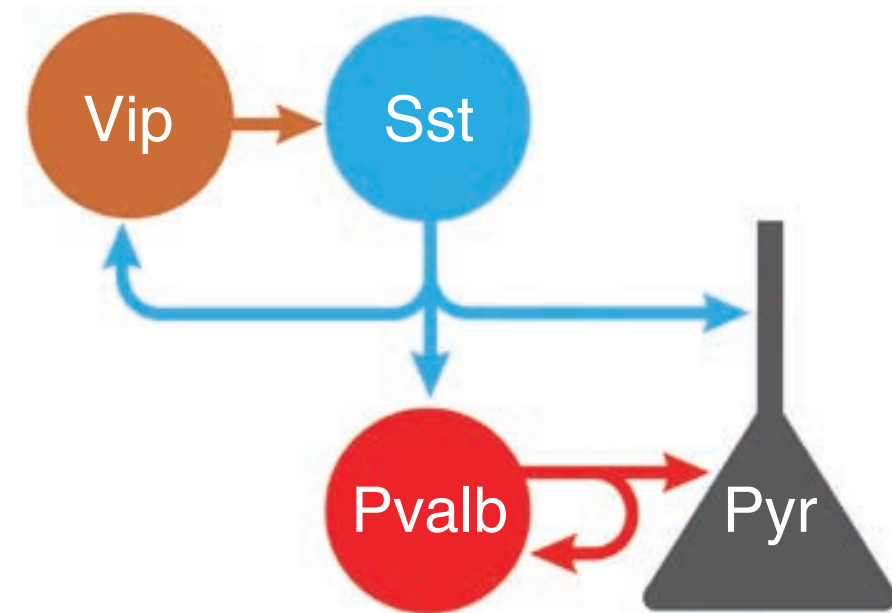
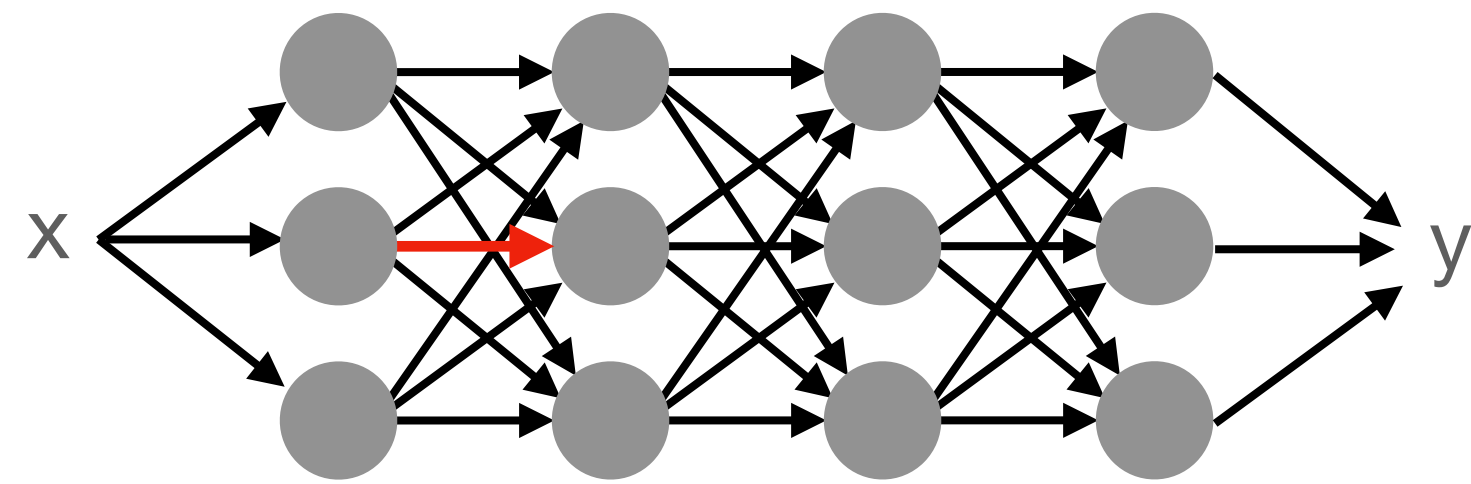
Burst-dependent synaptic plasticity can coordinate learning in hierarchical circuits

[Alexandre Payeur](#), [Jordan Guerguiev](#), [Friedemann Zenke](#), [Blake A. Richards](#)  & [Richard Naud](#) 

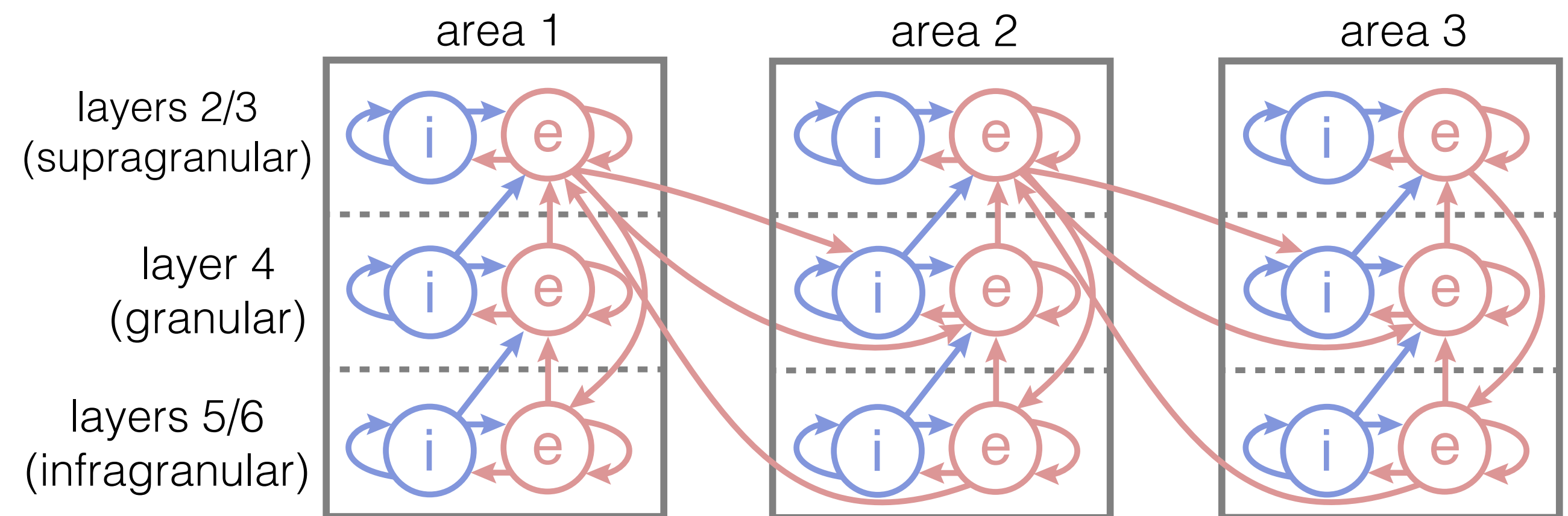
[Nature Neuroscience](#) **24**, 1010–1019 (2021) | [Cite this article](#)

Future plans

- Build a *framework* for meta-plasticity in different models with different tasks.

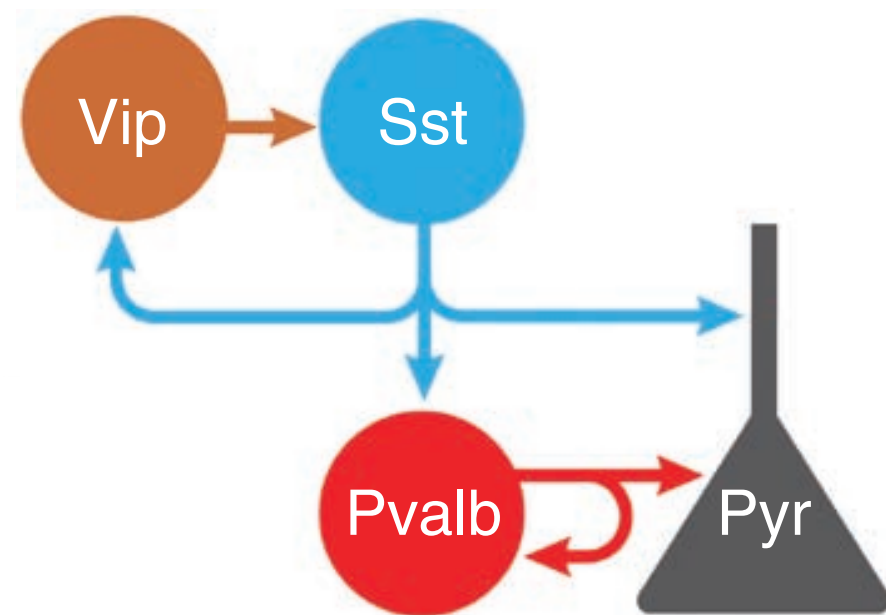
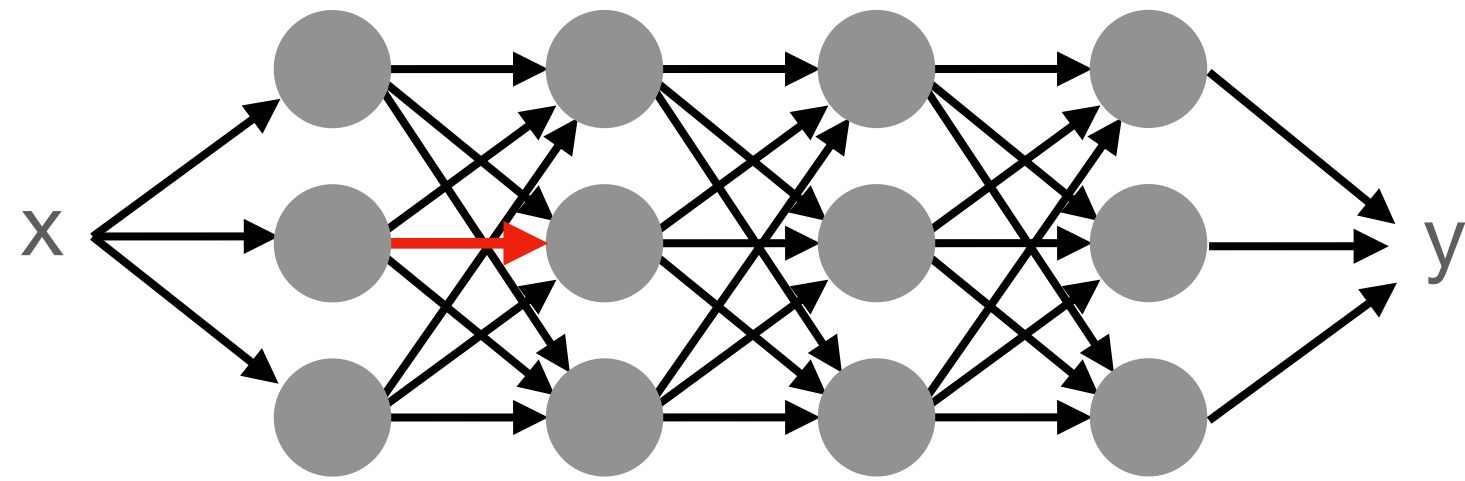


(Pfeffer, et al., 2013)



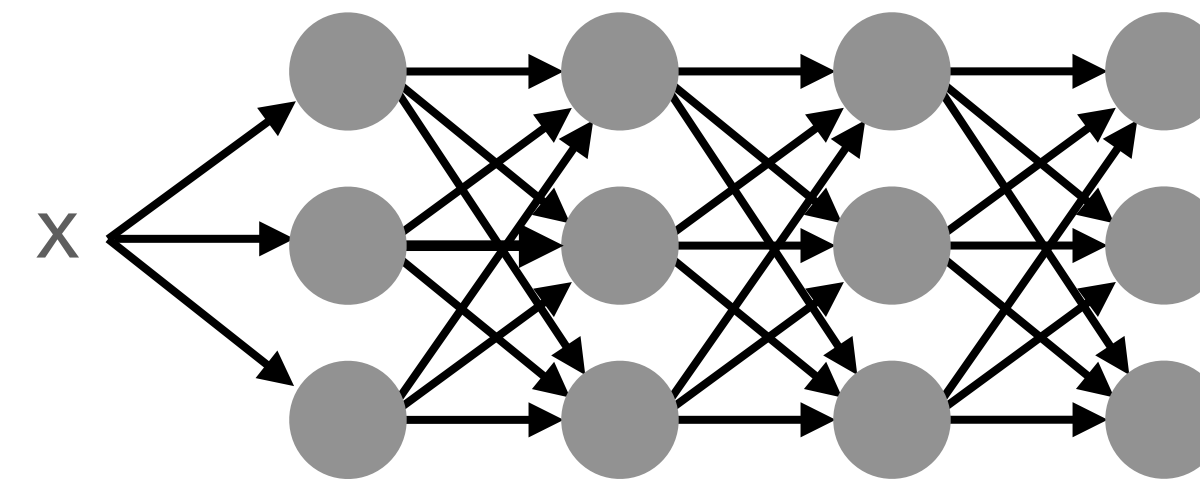
Future plans

- Build a *framework* for meta-plasticity in different models with different tasks.

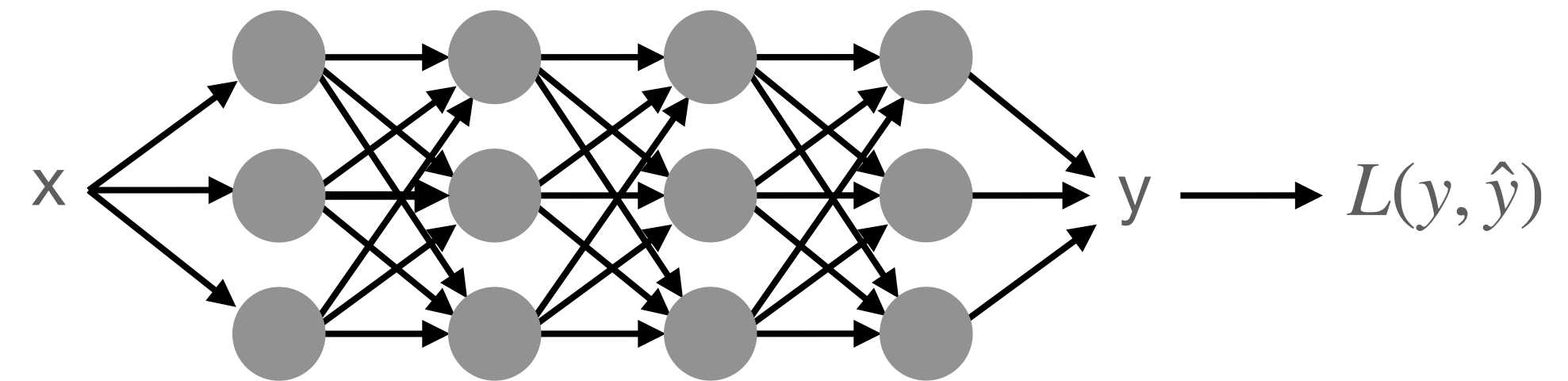


(Pfeffer, et al., 2013)

plasticity trains unsupervised representation



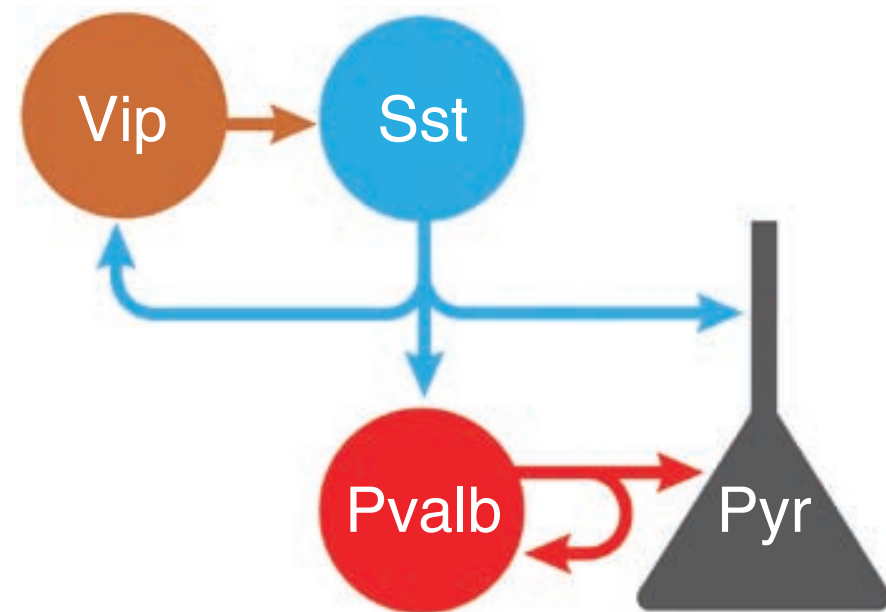
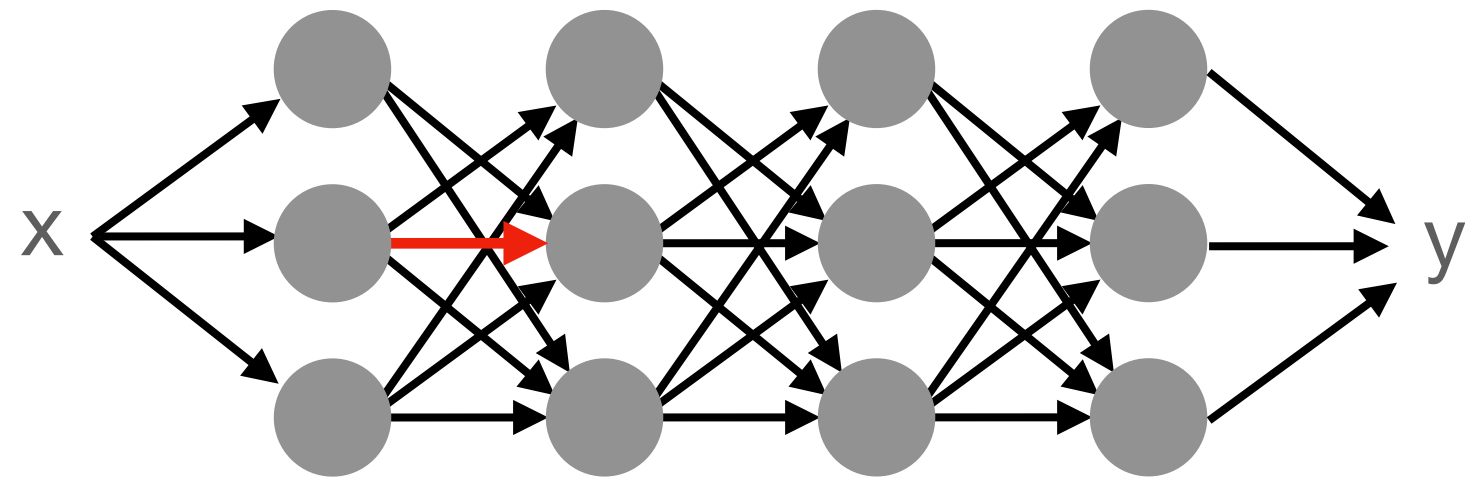
meta-loss evaluates suitability for supervised learning



Undergrad researcher: Sarah Duessing

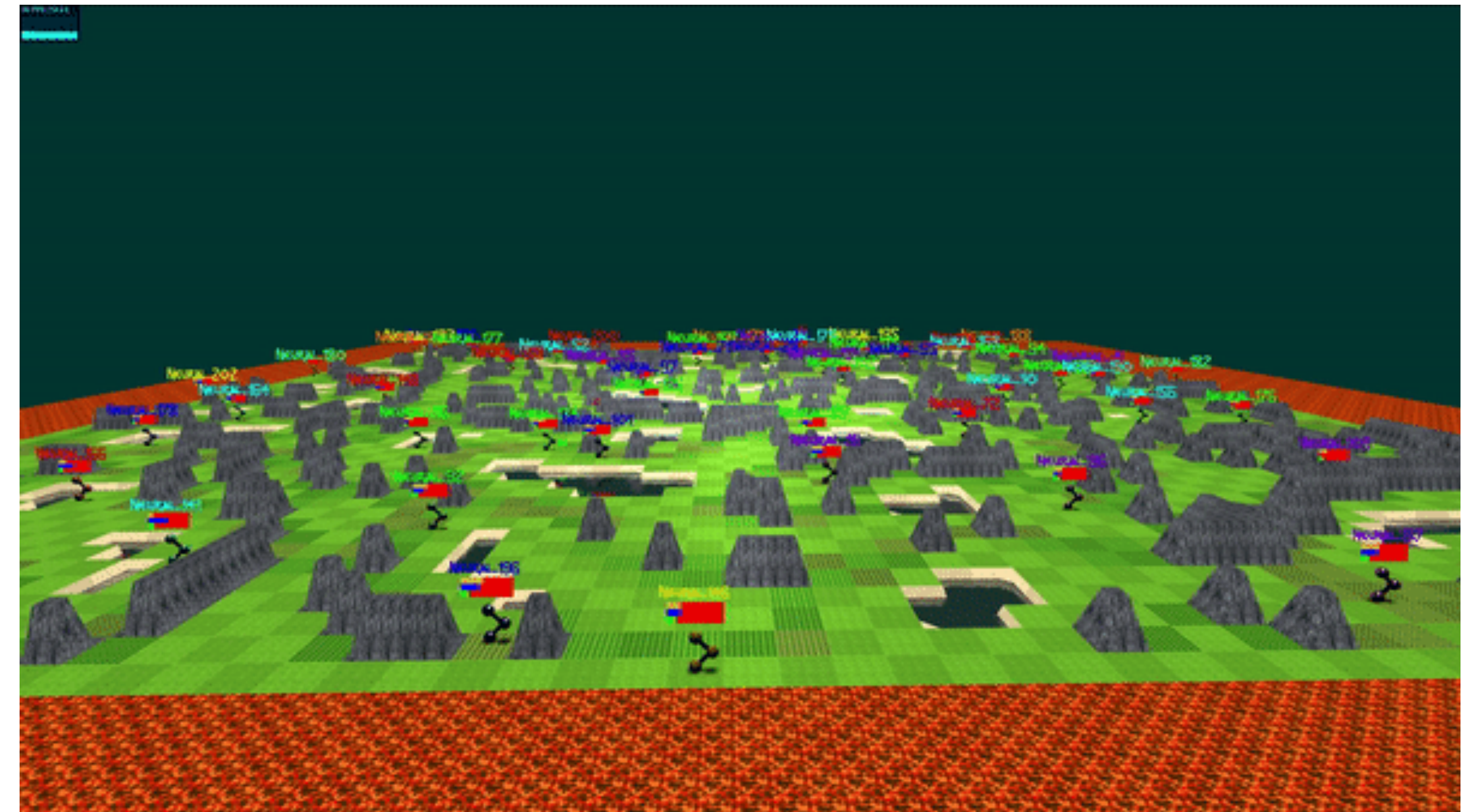
Future plans

- Build a *framework* for meta-plasticity in different models with different tasks.
- *Evolve* learning rules in multi-agent environments (gradient-free).



(Pfeffer, et al., 2013)

Neural MMO by OpenAI



Undergrad researcher: Sarah Duessing