

# Correlative Information Maximization: A Biologically Plausible Approach to Supervised Deep Neural Networks without Weight Symmetry

Neurips 2023

Bariscan Bozkurt<sup>1,2,3</sup>

Cengiz Pehlevan<sup>4,5</sup>

Alper T. Erdogan<sup>1,2</sup>

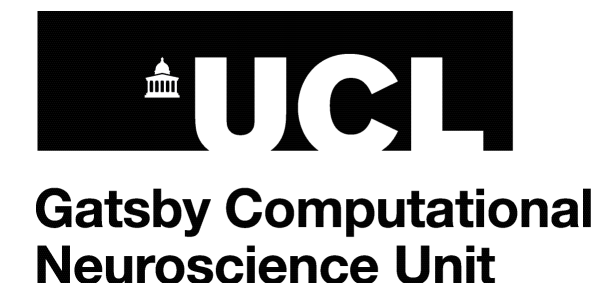
<sup>1</sup> KUIS AI Center, Koc University, Istanbul, Turkey

<sup>2</sup> Electrical and Electronics Engineering Department, Koc University, Istanbul, Turkey

<sup>3</sup> Gatsby Computational Neuroscience Unit, UCL, United Kingdom

<sup>4</sup> John A. Paulson School of Engineering & Applied Sciences and Center for Brain Science, Harvard University, Cambridge, MA 02138, USA

<sup>5</sup> Kempner Institute for the Study of Natural and Artificial Intelligence



# Goal

Introduce a **biologically plausible neural network framework** grounded on **information theory** offering,

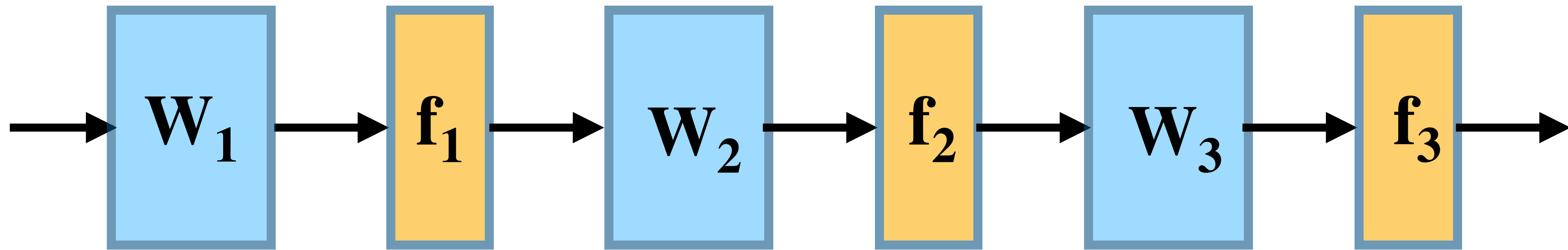
- ◆ A principled solution to the **weight symmetry** problem,
- ◆ A normative approach for deriving networks with **multi-compartment neurons**.

# Outline

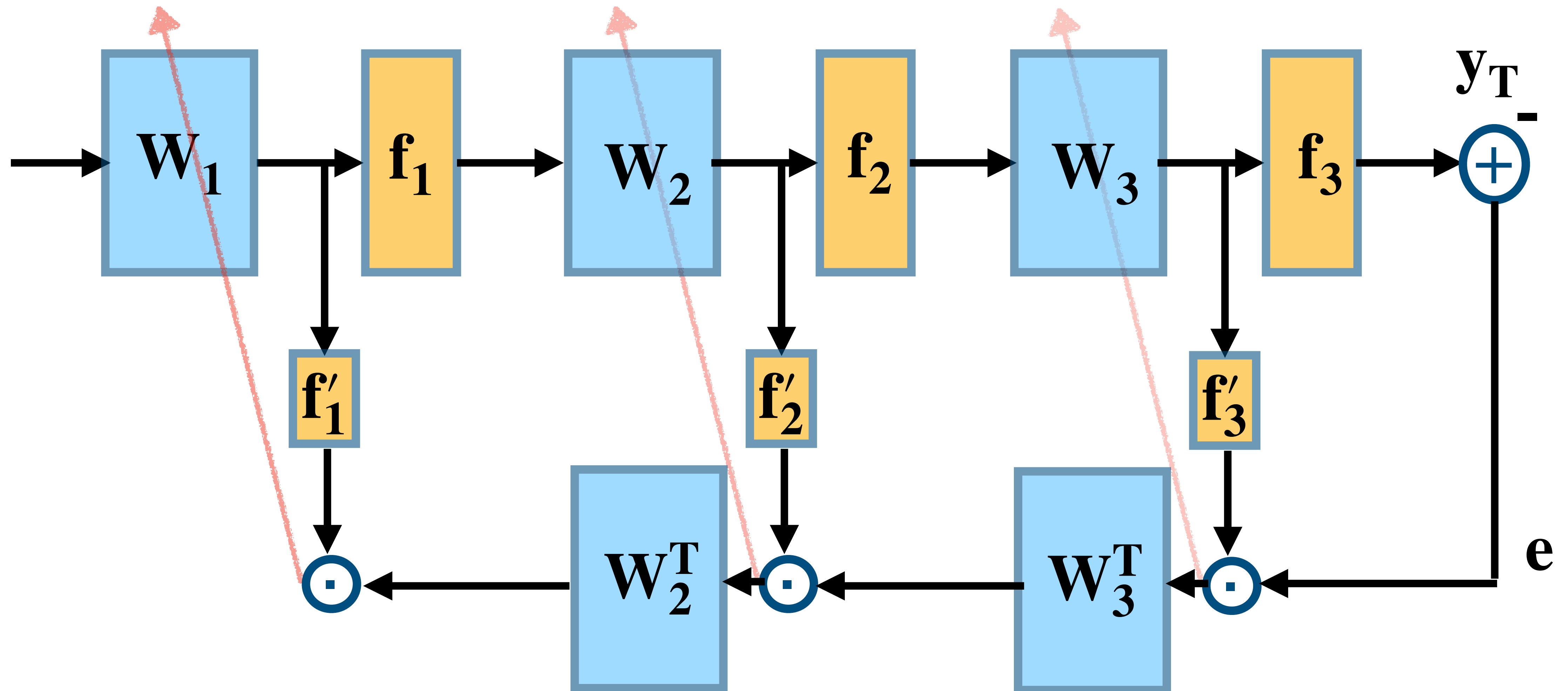
- ◆ Weight Symmetry (Transport) problem
- ◆ **Correlative Information Maximization** (CorInfoMax) criterion
- ◆ Network Data Model
- ◆ Network Structure and Dynamics
- ◆ Numerical Examples

# Backpropagation: Weight Symmetry Issue

A Multi-Layer Neural Network



# Backpropagation: Weight Symmetry Issue



# Correlative Information Maximization (CorInfoMax) Criterion

## Correlative Mutual Information (CMI)

**a, b** : Random vectors

$$I^{\rightarrow}(\epsilon_k)(\mathbf{a}, \mathbf{b}) = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{a}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{e}_{\mathbf{a}|\mathbf{b}}}[t] + \epsilon_k \mathbf{I})$$

The correlation matrix of **a**

The error of the best linear MMSE estimator of **a** from **b**

# Correlative Information Maximization (CorInfoMax) Criterion

## Correlative Mutual Information (CMI)

**a, b** : Random vectors

$$I^{(\epsilon_k)}(\mathbf{a}, \mathbf{b}) = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{a}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{e}_{\mathbf{a}|\mathbf{b}}}[t] + \epsilon_k \mathbf{I})$$

The correlation matrix of **a**

The error of the best linear MMSE estimator of **a** from **b**

Maximize  $I^{(\epsilon_k)}(\mathbf{a}, \mathbf{b}) \implies$  Maximize correlation (linear dependence) between **a** and **b**

# Network Data Model

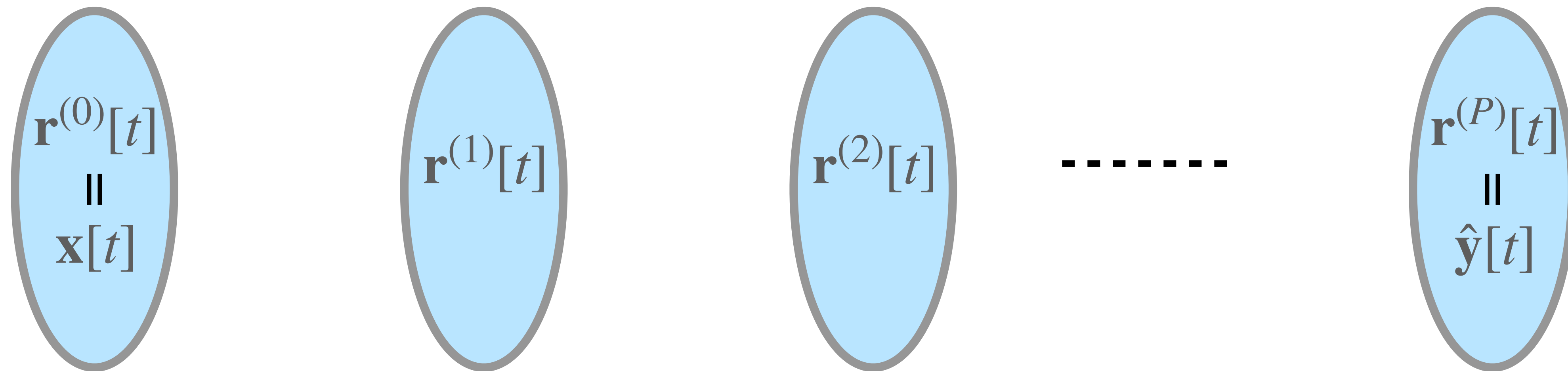
$\mathbf{x}[t]$  : input sample



# Network Data Model

$\mathbf{x}[t]$  : input sample

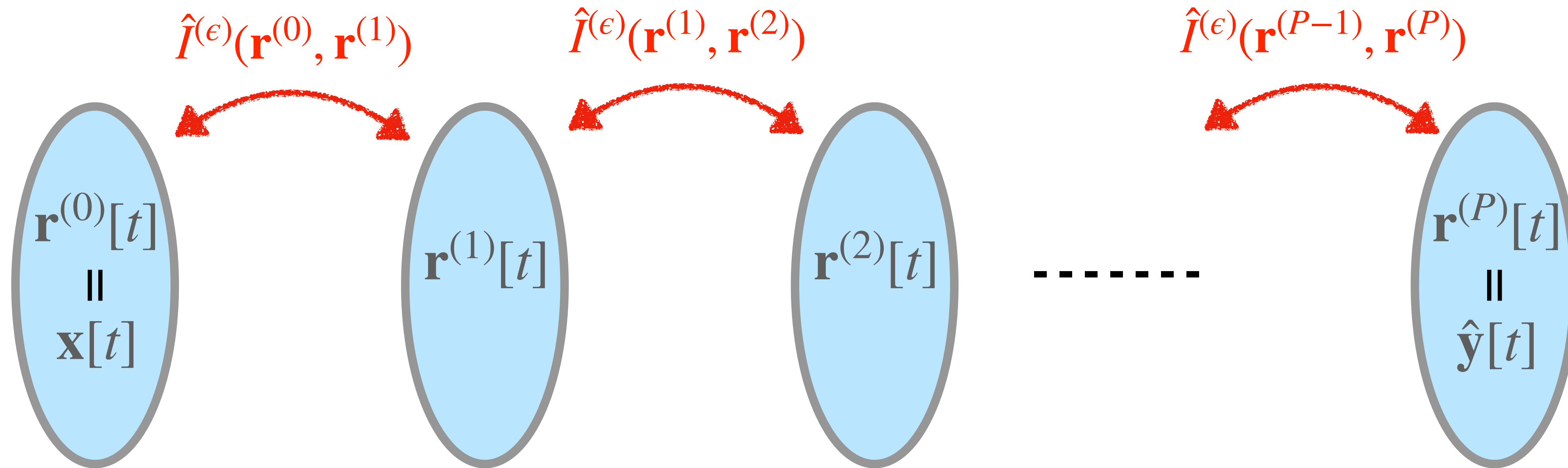
$\mathbf{r}^{(k)}[t]$  : layer activations



# Correlative Information Maximization Criterion

Objective:

$$\sum_{k=0}^{P-1} \hat{I}^{(\epsilon)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})$$



# Correlative Information Maximization Criterion

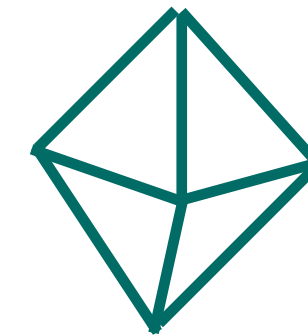
Objective:

$$\sum_{k=0}^{P-1} \hat{I}^{(\epsilon)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})$$

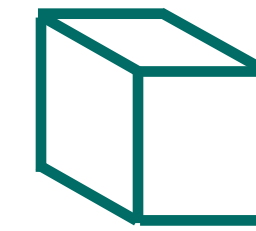
Constraints:

$$\mathbf{r}^{(k)} \in \mathcal{P}_k$$

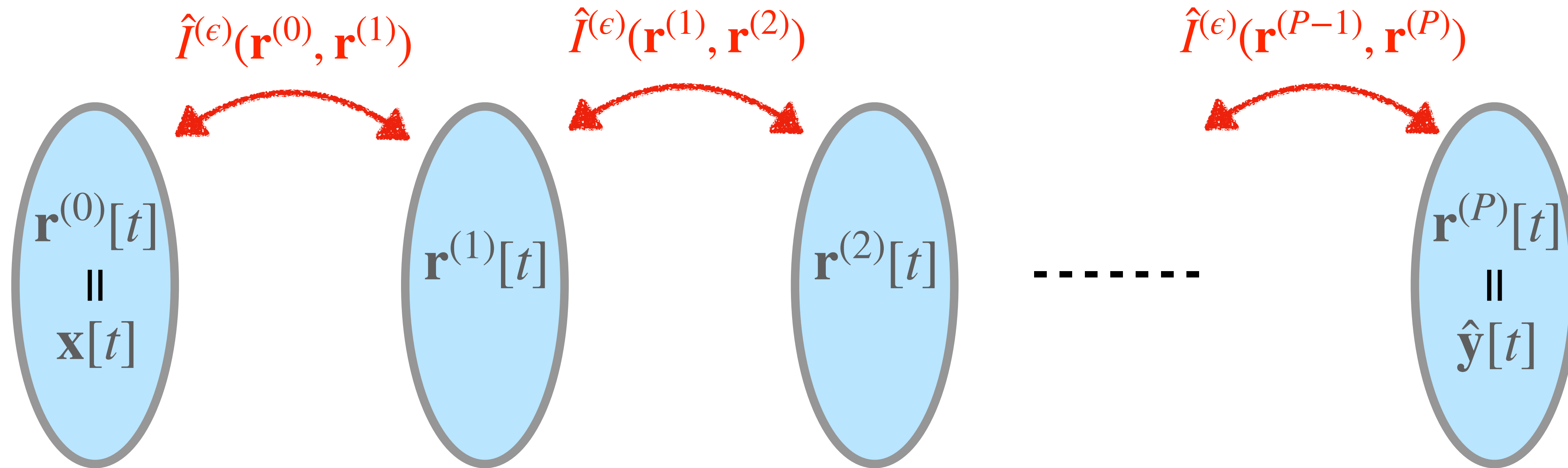
Polytope Examples



$\ell_1$  - norm ball  
Sparse



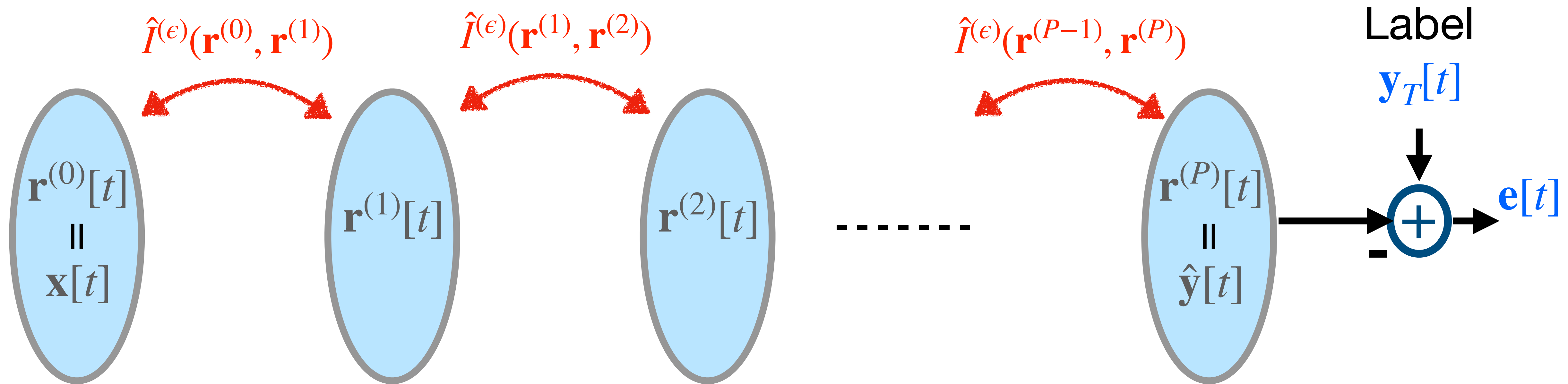
$\ell_\infty$  - norm ball  $\cap \mathbb{R}_+$   
Anti-sparse &  
Nonnegative



# Correlative Information Maximization Criterion

Objective: 
$$\sum_{k=0}^{P-1} \hat{I}^{(\epsilon)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})$$

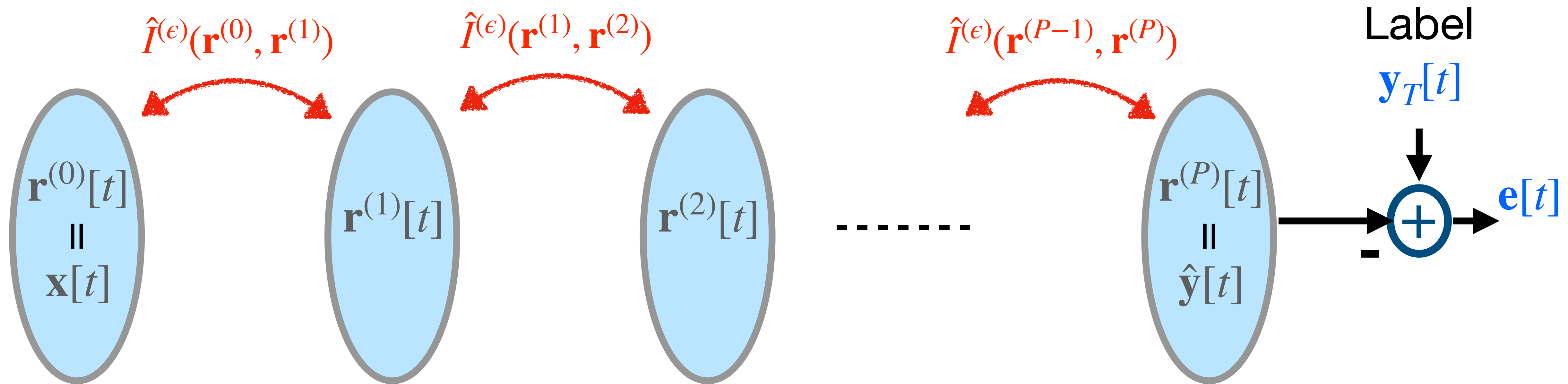
Constraints: 
$$\mathbf{r}^{(k)} \in \mathcal{P}_k$$



# Correlative Information Maximization Criterion

Objective: 
$$\sum_{k=0}^{P-1} \hat{I}^{(\epsilon)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)}) - \frac{\beta}{2} \|\mathbf{y}_T[t] - \mathbf{r}^{(P)}[t]\|_2^2$$

Constraints: 
$$\mathbf{r}^{(k)} \in \mathcal{P}_k$$



# Correlative Information Maximization Criterion

Two Alternative Expressions for the Correlative Mutual Information (CMI):

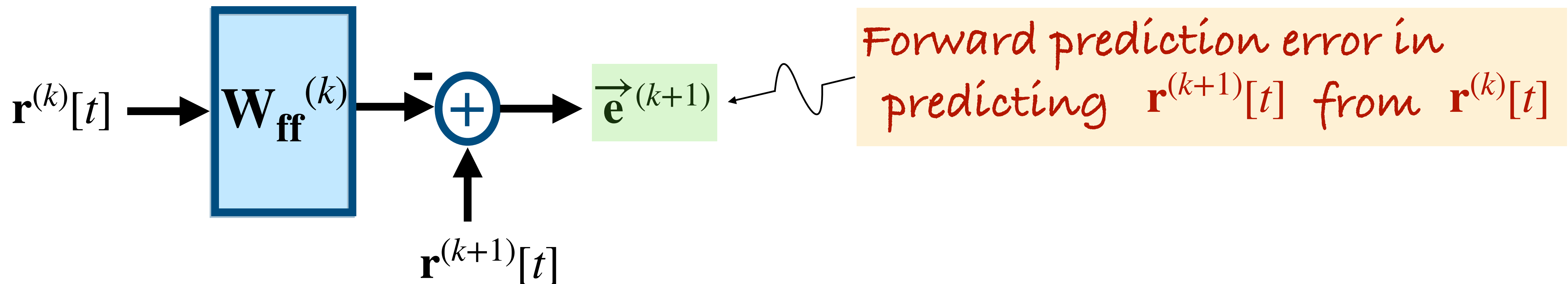
# Correlative Information Maximization Criterion

Two Alternative Expressions for the Correlative Mutual Information (CMI):

Alternative 1:

$$\vec{\hat{I}}^{(\epsilon_k)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t] = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{r}^{(k+1)}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\vec{\mathbf{e}}_*^{(k+1)}}[t] + \epsilon_k \mathbf{I})$$

where



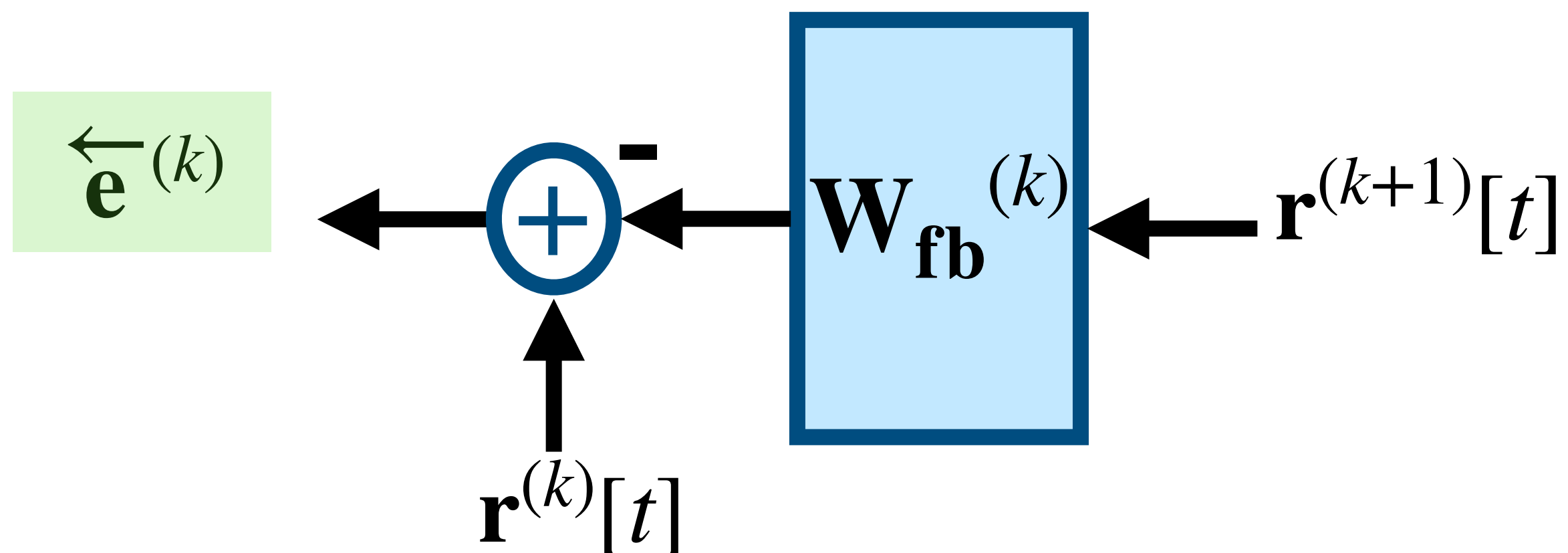
# Correlative Information Maximization Criterion

Two Alternative Expressions for the Correlative Mutual Information (CMI):

Alternative 2:

$$\hat{I}^{\leftarrow}(\epsilon_k)(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t] = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{r}^{(k)}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{e}_*^{\leftarrow}(k)}[t] + \epsilon_k \mathbf{I})$$

where



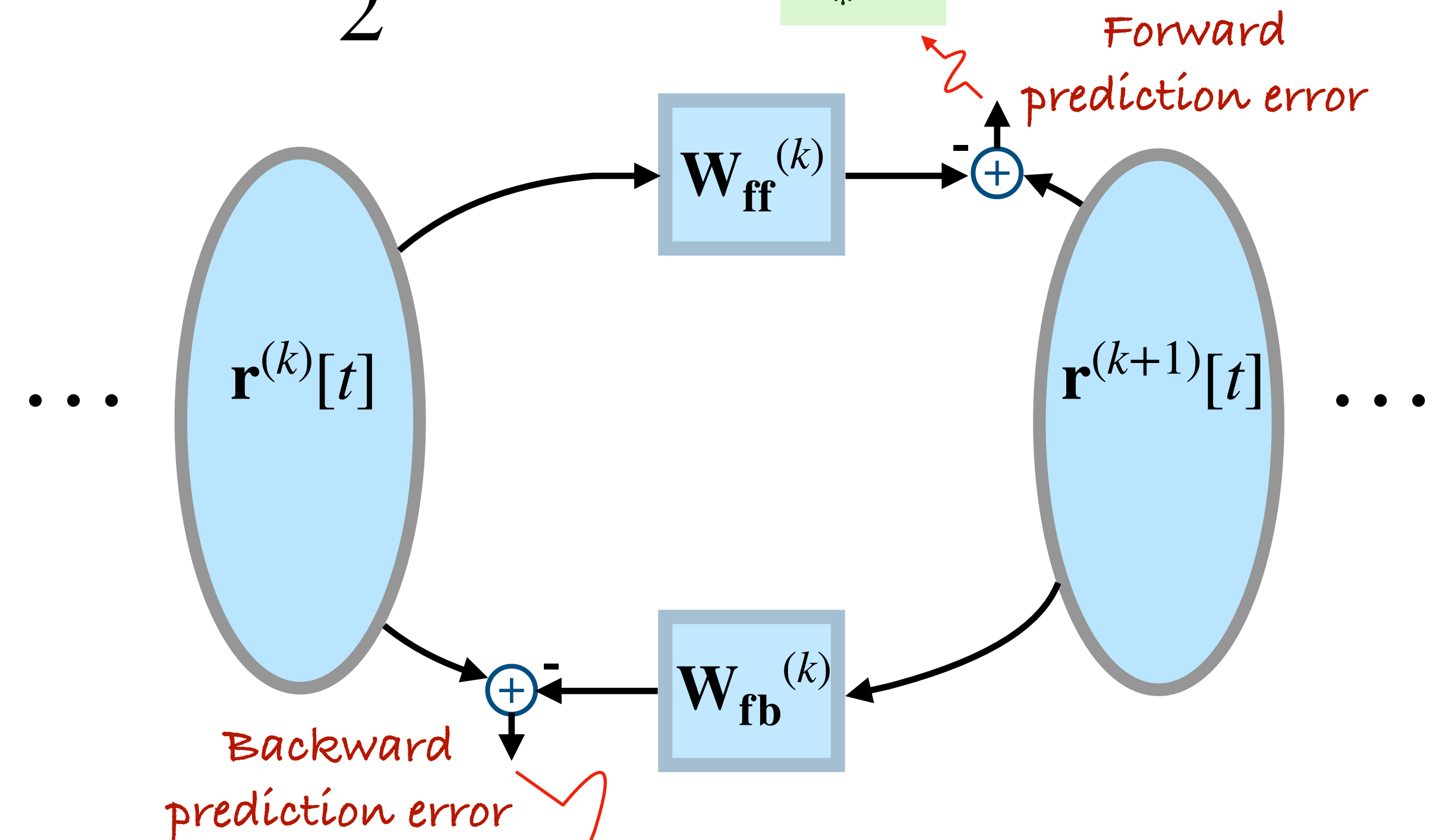
Backward prediction error in predicting  $\mathbf{r}^{(k)}[t]$  from  $\mathbf{r}^{(k+1)}[t]$



# Correlative Information Maximization Criterion

Two Alternative Expressions for the Correlative Mutual Information (CMI):

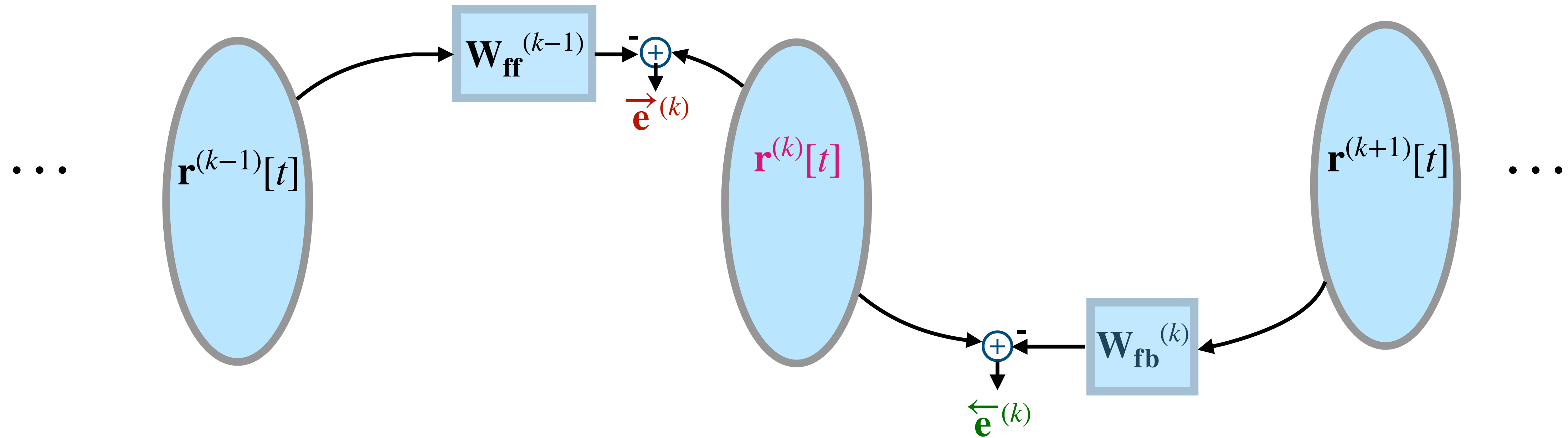
$$\vec{\hat{I}}^{(\epsilon_k)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t] = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{r}^{(k+1)}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\vec{\mathbf{e}}_*^{(k+1)}}[t] + \epsilon_k \mathbf{I})$$



$$\leftarrow \hat{I}^{(\epsilon_k)}(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t] = \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\mathbf{r}^{(k)}}[t] + \epsilon_k \mathbf{I}) - \frac{1}{2} \log \det(\hat{\mathbf{R}}_{\leftarrow \mathbf{e}}^{(k)}[t] + \epsilon_k \mathbf{I})$$

# Correlative Information Maximization Criterion

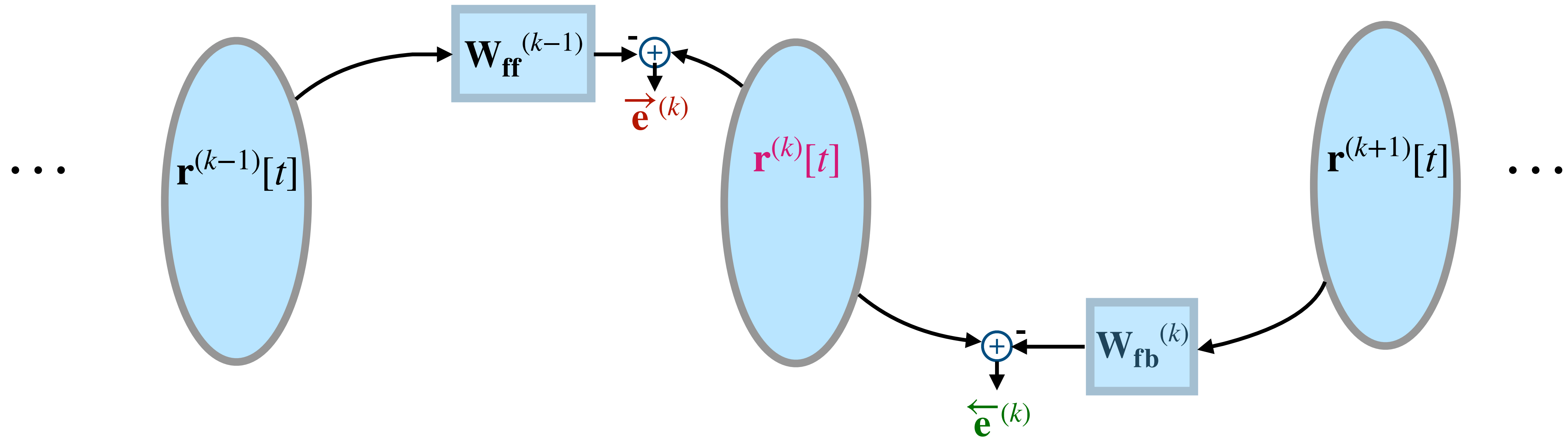
Maximization of CMI: Gradient wrt  $\mathbf{r}^{(k)}[t]$



$$\hat{J}_k(\mathbf{r}^{(k)}[t]) = \hat{I}^{\rightarrow}(\epsilon_{k-1})(\mathbf{r}^{(k-1)}, \mathbf{r}^{(k)})[t] + \hat{I}^{\leftarrow}(\epsilon_k)(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t]$$

# Correlative Information Maximization Criterion

Maximization of CMI: Gradient wrt  $\mathbf{r}^{(k)}[t]$



$$\nabla_{\mathbf{r}^{(k)}} \hat{J}_k(\mathbf{r}^{(k)}[t]) = \nabla_{\mathbf{r}^{(k)}} \hat{I}^{\rightarrow}(\epsilon_{k-1})(\mathbf{r}^{(k-1)}, \mathbf{r}^{(k)})[t] + \nabla_{\mathbf{r}^{(k)}} \hat{I}^{\leftarrow}(\epsilon_k)(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)})[t]$$

# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k}(\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k}(\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k} (\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k} (\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

$\mathbf{u}^{(k)}[t, s]$  : Somatic compartment membrane potential

# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k} (\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k} (\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

$\mathbf{u}^{(k)}[t, s]$  : Somatic compartment membrane potential

$\mathbf{v}_A^{(k)}[t, s] = \mathbf{M}^{(k)}[t] \mathbf{r}^{(k)}[t; s] + \mathbf{W}_{fb}^{(k)}[t] \mathbf{r}^{(k+1)}[t; s]$  Apical compartment membrane potential

# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k} (\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k} (\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

$\mathbf{u}^{(k)}[t, s]$  : Somatic compartment membrane potential

$\mathbf{v}_B^{(k)}[t, s] = \mathbf{W}_{ff}^{(k-1)}[t] \mathbf{r}^{(k-1)}[t; s]$  Basal compartment membrane potential

$\mathbf{v}_A^{(k)}[t, s] = \mathbf{M}^{(k)}[t] \mathbf{r}^{(k)}[t; s] + \mathbf{W}_{fb}^{(k)}[t] \mathbf{r}^{(k+1)}[t; s]$  Apical compartment membrane potential

# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k} (\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k} (\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

$$g_{A,k} = \frac{1}{\epsilon_{k-1}} \quad \begin{array}{l} \text{Apical-Soma} \\ \text{Conductance} \end{array}$$

$$g_{B,k} = \frac{1}{\epsilon_k} \quad \begin{array}{l} \text{Basal-Soma} \\ \text{Conductance} \end{array}$$

$$g_{lk} : \quad \begin{array}{l} \text{Leakage} \\ \text{Conductance} \end{array}$$



# Network Structure and Neural Dynamics

## Multi-compartmental neural network

Projected gradient ascent based neural-dynamics:

Sample index

Continuous time index

$$\tau_{\mathbf{u}} \frac{d\mathbf{u}^{(k)}[t; s]}{ds} = -g_{lk} \mathbf{u}^{(k)}[t; s] + g_{A,k} (\mathbf{v}_A^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s]) + g_{B,k} (\mathbf{v}_B^{(k)}[t; s] - \mathbf{u}^{(k)}[t; s])$$

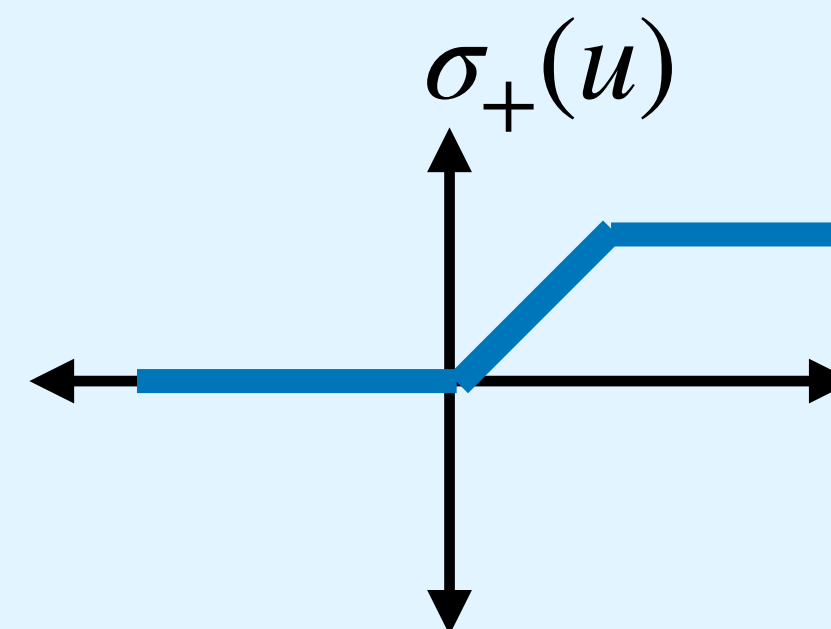
$$\mathbf{r}^{(k)}[t; s] = \sigma_+(\mathbf{u}^{(k)}[t; s])$$

where

$$g_{A,k} = \frac{1}{\epsilon_{k-1}} \quad \begin{array}{l} \text{Apical-Soma} \\ \text{Conductance} \end{array}$$

$$g_{B,k} = \frac{1}{\epsilon_k} \quad \begin{array}{l} \text{Basal-Soma} \\ \text{Conductance} \end{array}$$

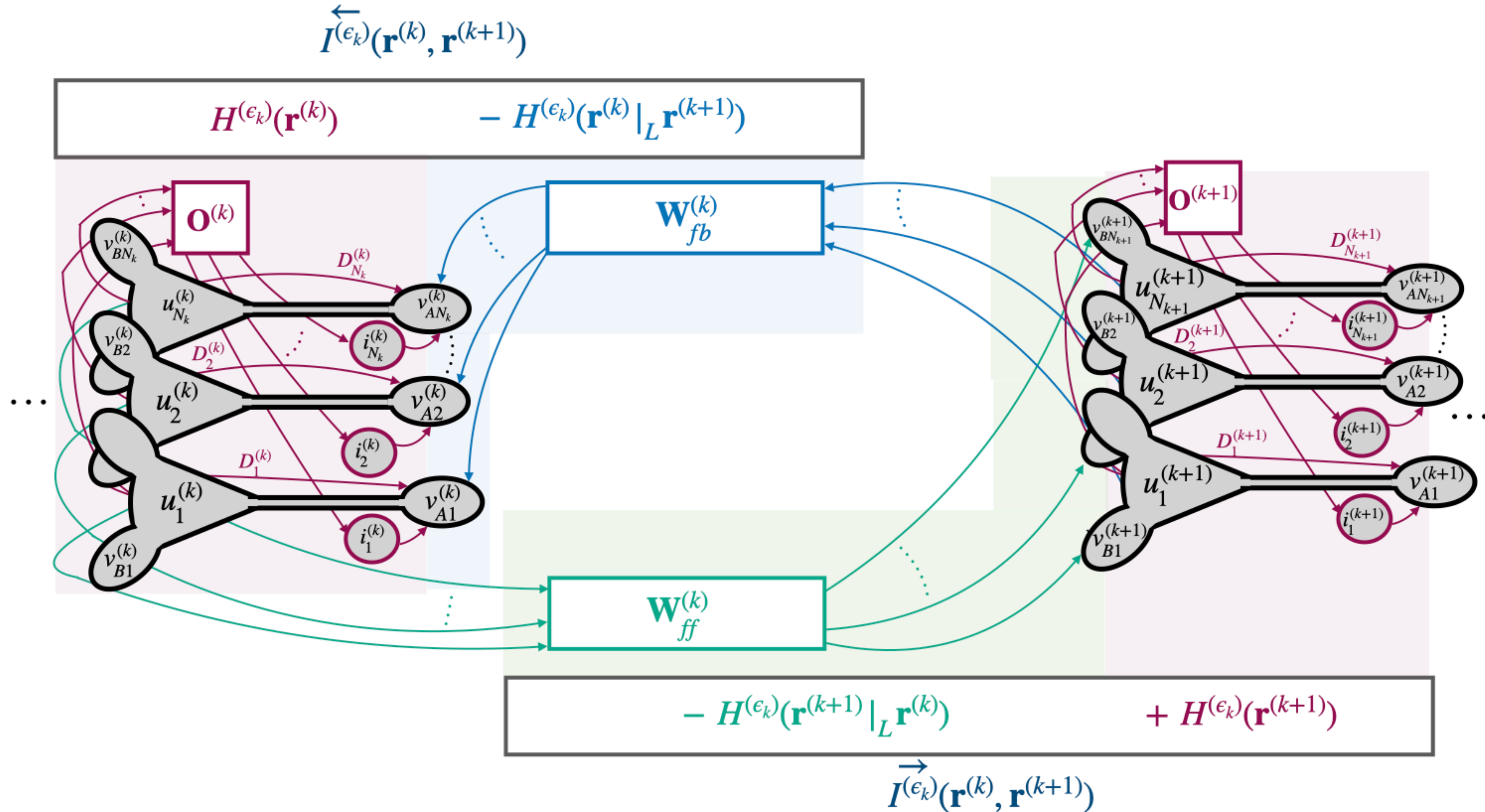
$$g_{lk} : \quad \begin{array}{l} \text{Leakage} \\ \text{Conductance} \end{array}$$



Projection  
onto  
Hypercube

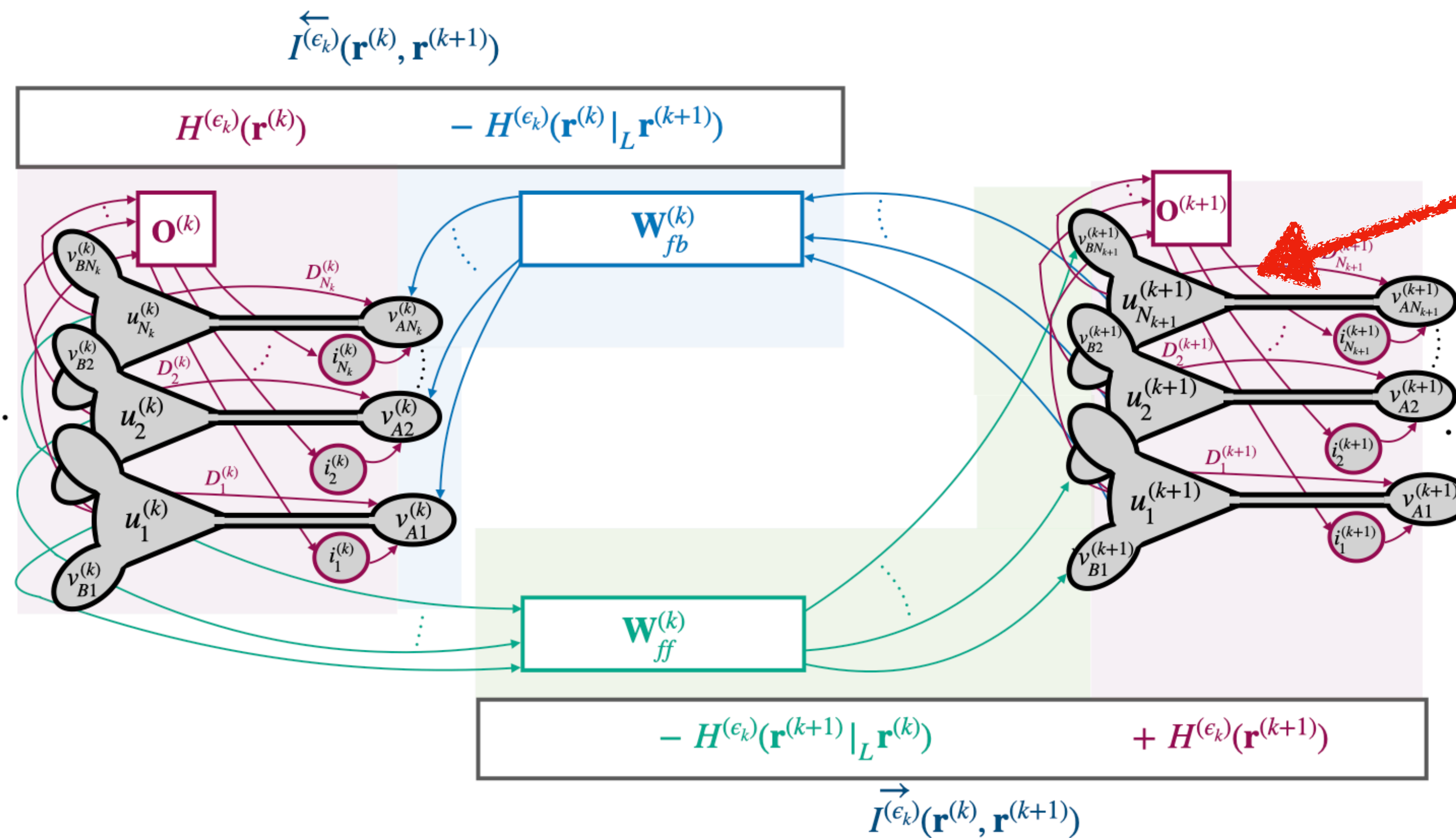
# Network Structure and Neural Dynamics

## Multi-compartmental neural network



# Network Structure and Neural Dynamics

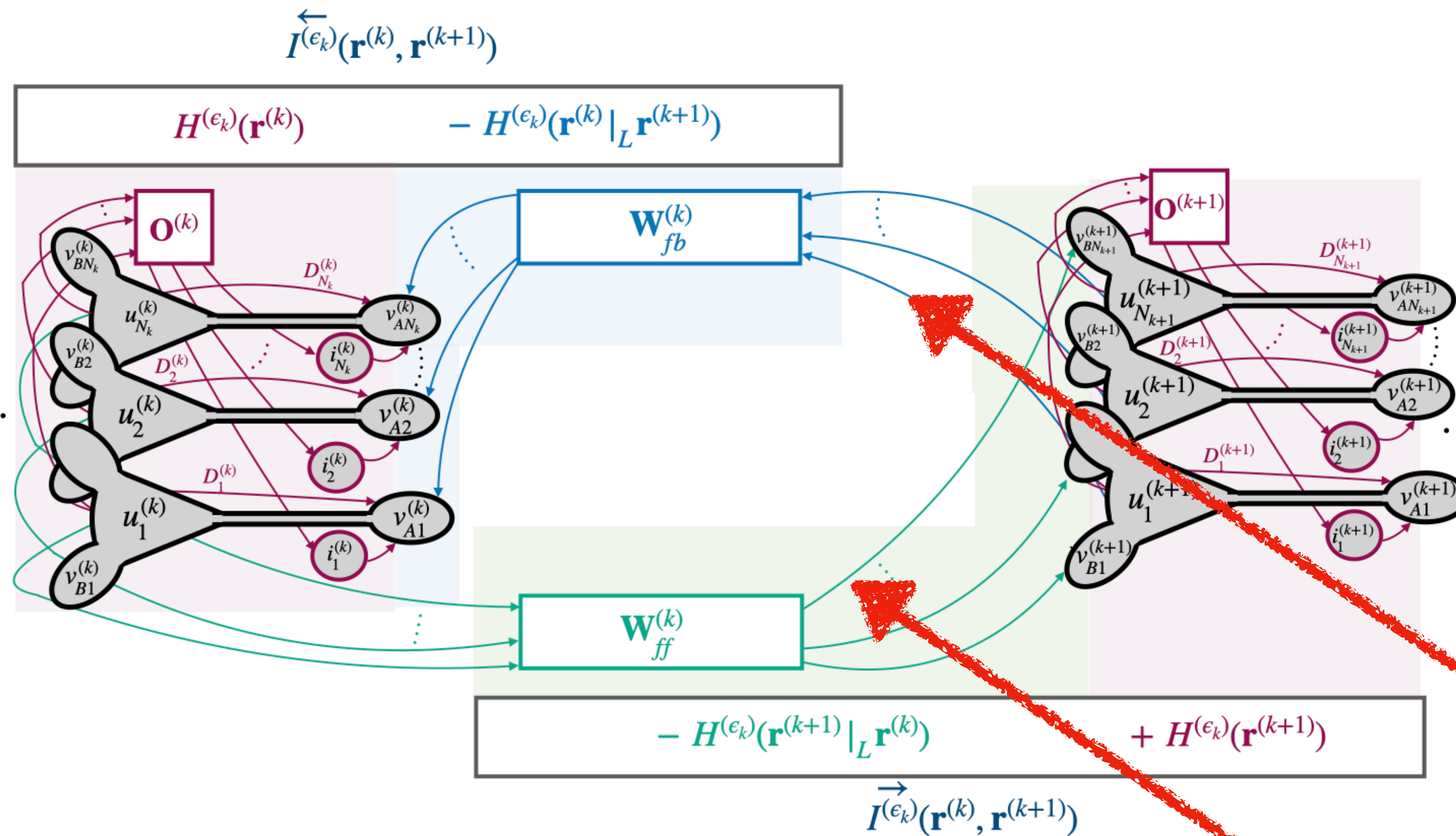
## Multi-compartmental neural network



Lateral connections  
to maximize layer entropy:  
utilization of representation  
dimensions, avoid degeneracy

# Network Structure and Neural Dynamics

## Multi-compartmental neural network



Feedforward/Feedback connections to minimize conditional entropy: Facilitate bidirectional information flow, reduce redundancy

# Network Structure and Neural Dynamics

## Learning Dynamics

$$\delta \mathbf{W}_{ff}^{(k)}[t] \propto \frac{1}{\beta'} \left( \left( \vec{\mathbf{e}}^{(k+1)}[t] \mathbf{r}^{(k)}[t]^T \right) \Big|_{\beta=\beta'} - \left( \vec{\mathbf{e}}^{(k+1)}[t] \mathbf{r}^{(k)}[t]^T \right) \Big|_{\beta=0} \right)$$

$$\delta \mathbf{W}_{fb}^{(k)}[t] \propto \frac{1}{\beta'} \left( \left( \overleftarrow{\mathbf{e}}^{(k)}[t] \mathbf{r}^{(k+1)}[t]^T \right) \Big|_{\beta=\beta'} - \left( \overleftarrow{\mathbf{e}}^{(k)}[t] \mathbf{r}^{(k+1)}[t]^T \right) \Big|_{\beta=0} \right)$$

$$\mathbf{B}^{(k)}[t+1] = \lambda_r^{-1} (\mathbf{B}^{(k)}[t] - \gamma \mathbf{z}^{(k)}[t] \mathbf{z}^{(k)}[t]^T), \text{ where } \mathbf{z}^{(k)} = \mathbf{B}^{(k)}[t] \mathbf{r}^{(k)}[t] \Big|_{\beta=\beta'}$$

# Numerical Examples

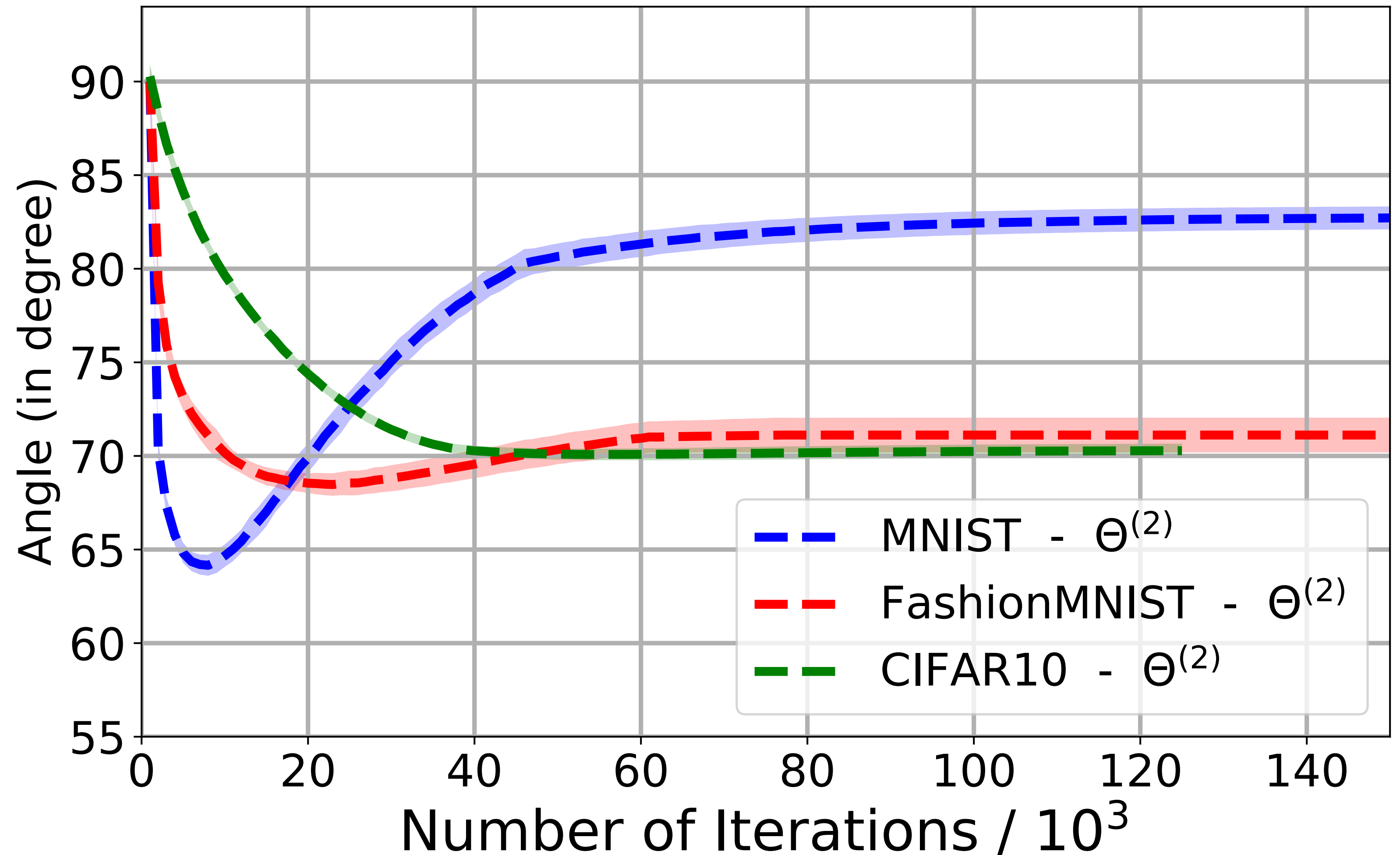
- **MNIST** and **Fashion MNIST** Datasets: Layer sizes of 784, 500,10
- **CIFAR 10** Dataset: Layer sizes of 3072,1000,10

	MNIST	FashionMNIST	CIFAR10
<b>CorInfoMax-<math>\mathcal{B}_{\infty,+}</math></b>	97.62 $\pm$ 0.1	88.14 $\pm$ 0.3	51.86 $\pm$ 0.3
<b>CorInfoMax-<math>\mathcal{B}_{1,+}</math></b>	97.71 $\pm$ 0.1	88.09 $\pm$ 0.1	51.19 $\pm$ 0.4
EP	97.61 $\pm$ 0.1	88.06 $\pm$ 0.7	49.28 $\pm$ 0.5
CSM	98.08 $\pm$ 0.1	88.73 $\pm$ 0.2	40.79*
PC	98.17 $\pm$ 0.2	89.31 $\pm$ 0.4	-
PC-Nudge	97.71 $\pm$ 0.1	88.49 $\pm$ 0.3	48.58 $\pm$ 0.7
Feedback Alignment (with MSE Loss)	97.99 $\pm$ 0.03	88.72 $\pm$ 0.5	50.75 $\pm$ 0.4
Feedback Alignment (with CrossEntropy Loss)	97.95 $\pm$ 0.08	88.38 $\pm$ 0.9	52.37 $\pm$ 0.4
BP (with MSE Loss)	97.58 $\pm$ 0.01	88.39 $\pm$ 0.1	52.75 $\pm$ 0.1
BP (with CrossEntropy Loss)	98.27 $\pm$ 0.03	89.41 $\pm$ 0.2	53.96 $\pm$ 0.3

# Numerical Examples

**Confirming Asymmetry:** Angle between feedforward and the transpose of the feedback weights

$$\Theta^{(k)} = \arccos \left( \frac{\text{Tr} \left( \mathbf{W}_{ff}^{(k)} \mathbf{W}_{fb}^{(k)} \right)}{\|\mathbf{W}_{ff}^{(k)}\|_F \|\mathbf{W}_{fb}^{(k)}\|_F} \right)$$



# Conclusions

**CorInfoMax** is offered as an information theory-based principled framework:

- Networks of segregated neurons with recurrent and asymmetric feedback/feedforward connections governed by local learning rules naturally emerge,
- Useful in obtaining potential insights such as
  - ◆ the role of lateral connections in embedding space expansion and avoiding degeneracy,
  - ◆ feedback and feedforward connections for prediction to reduce redundancy,
  - ◆ activation functions/interneurons to shape feature space and compress.



**Thank You!**