

# Optimal Learning Protocols via Statistical Physics and Control Theory

Francesco Mori



UNIVERSITY OF  
OXFORD

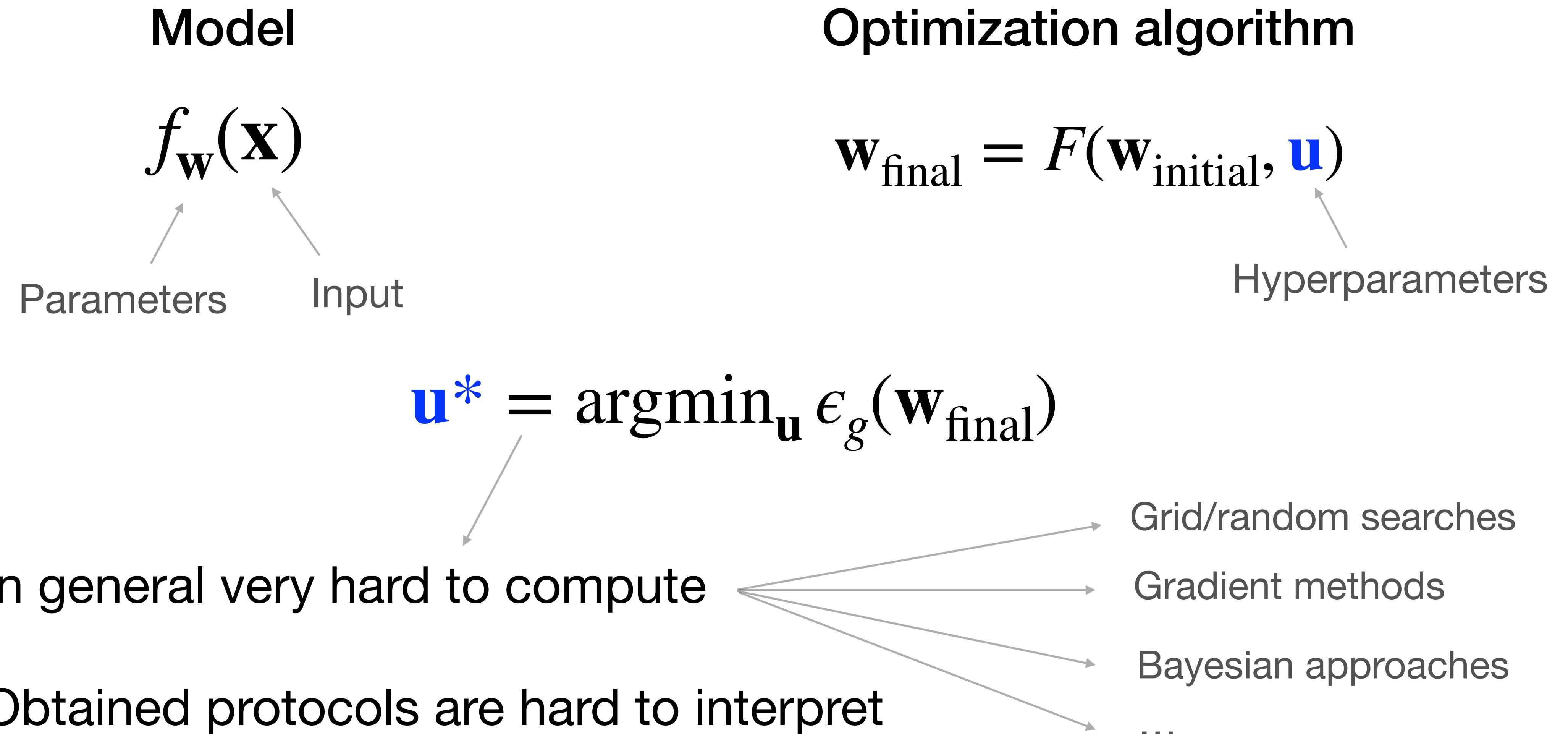


NEW COLLEGE  
OXFORD

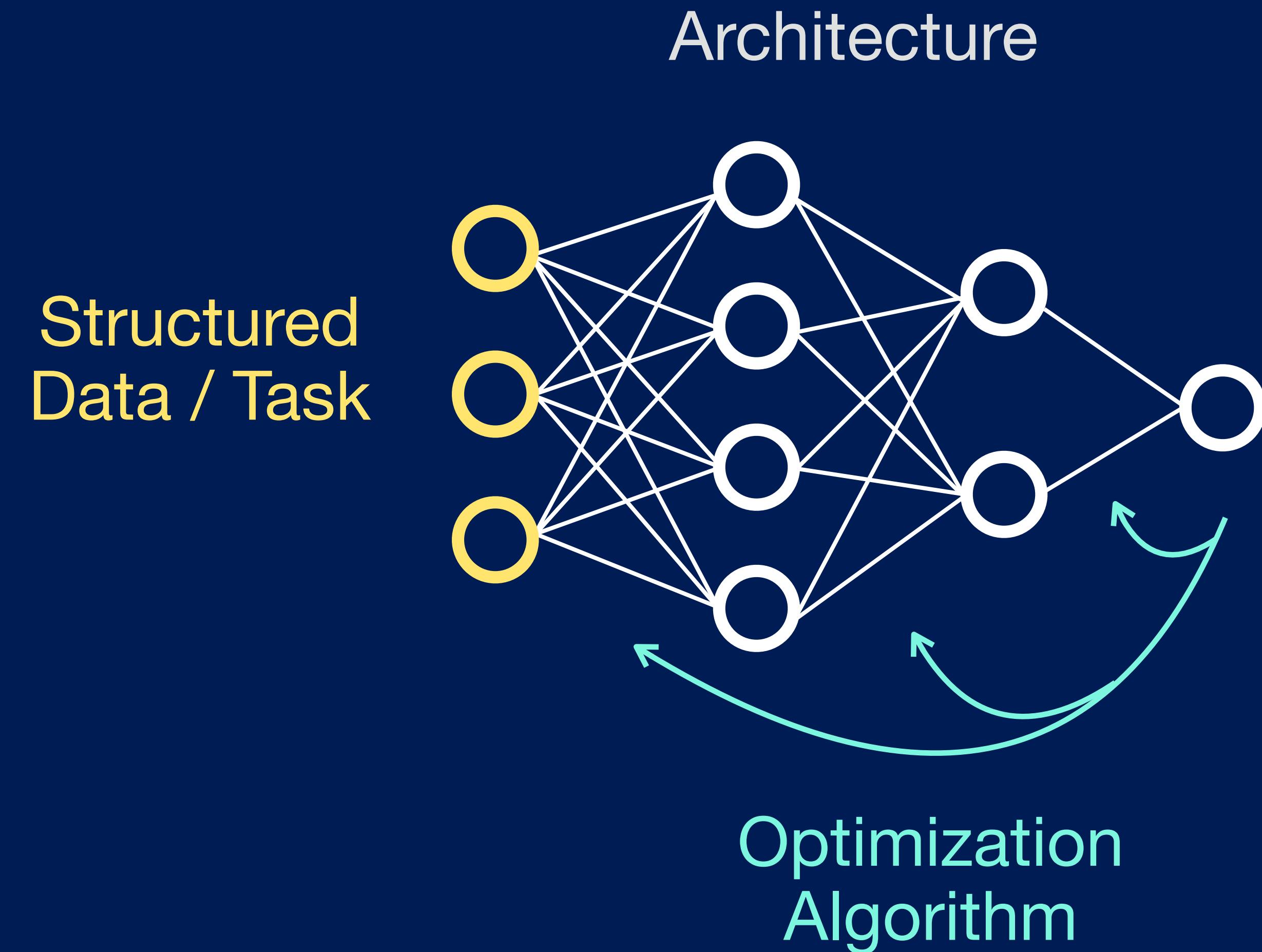
*Statistical Physics & Machine Learning: moving forward*

Institut d'Études Scientifiques de Cargèse

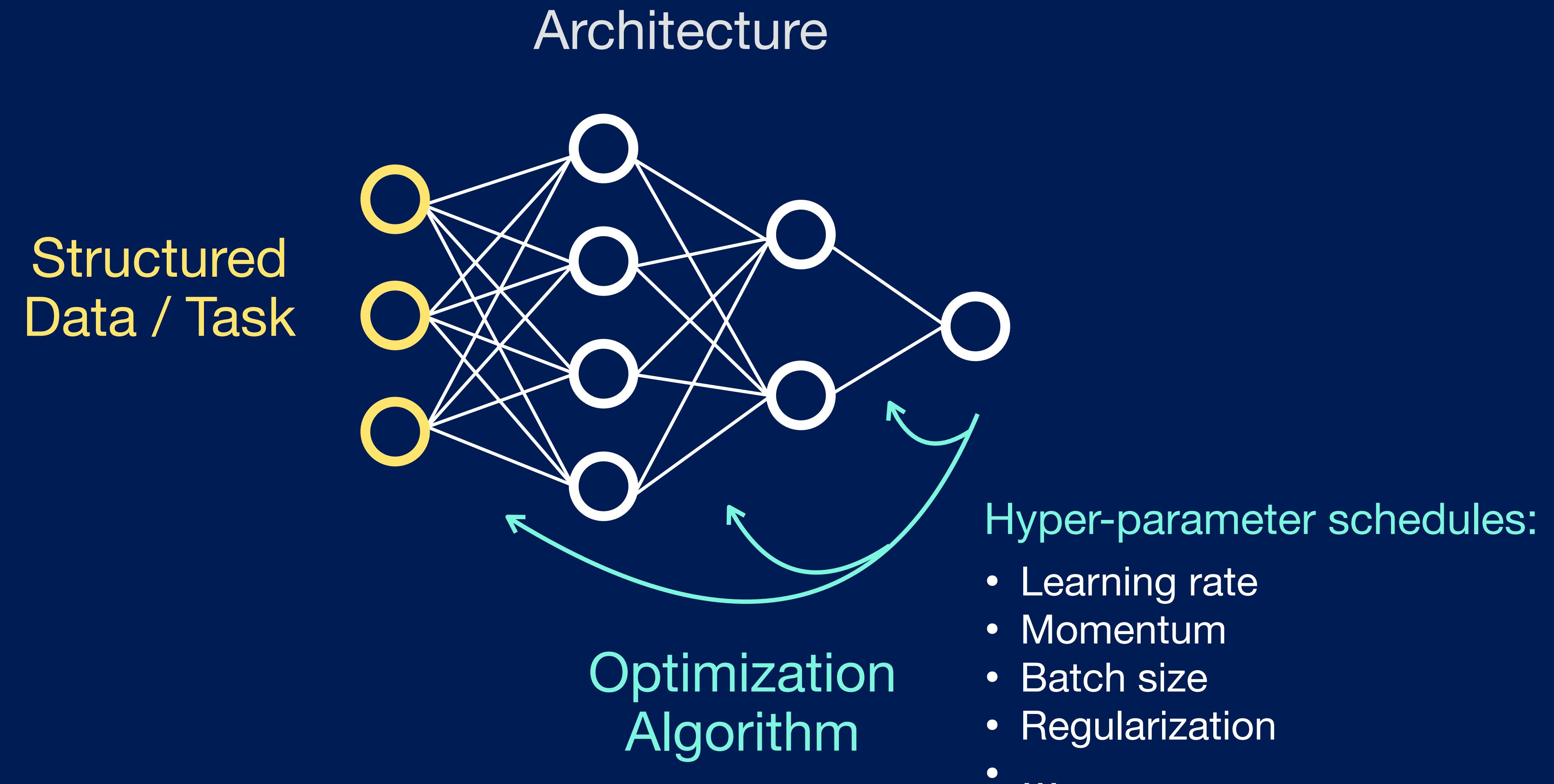
# General setup - Hyperparameter Optimization / Meta Learning



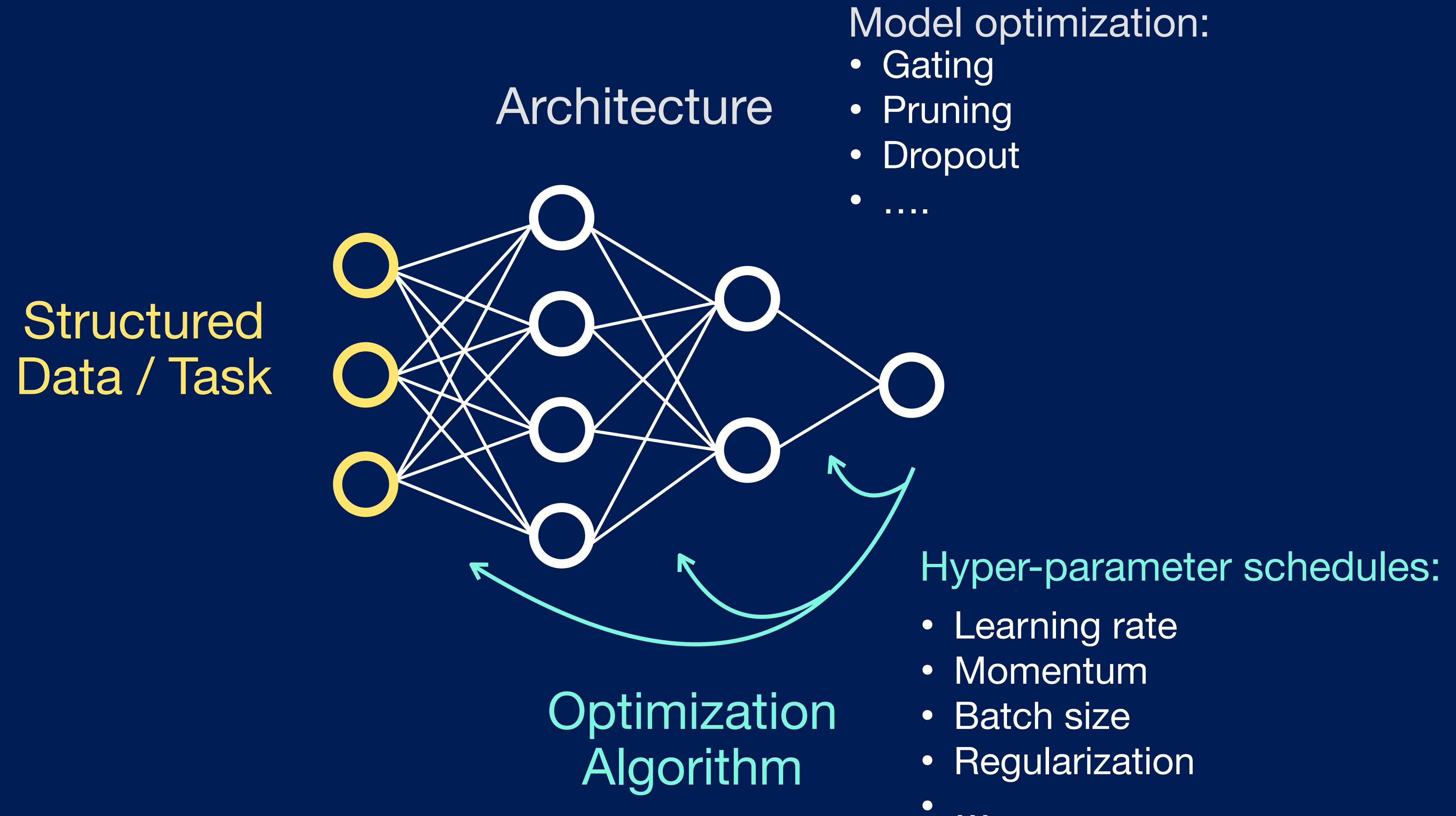
# Learning protocols for neural networks



# Learning protocols for neural networks



# Learning protocols for neural networks

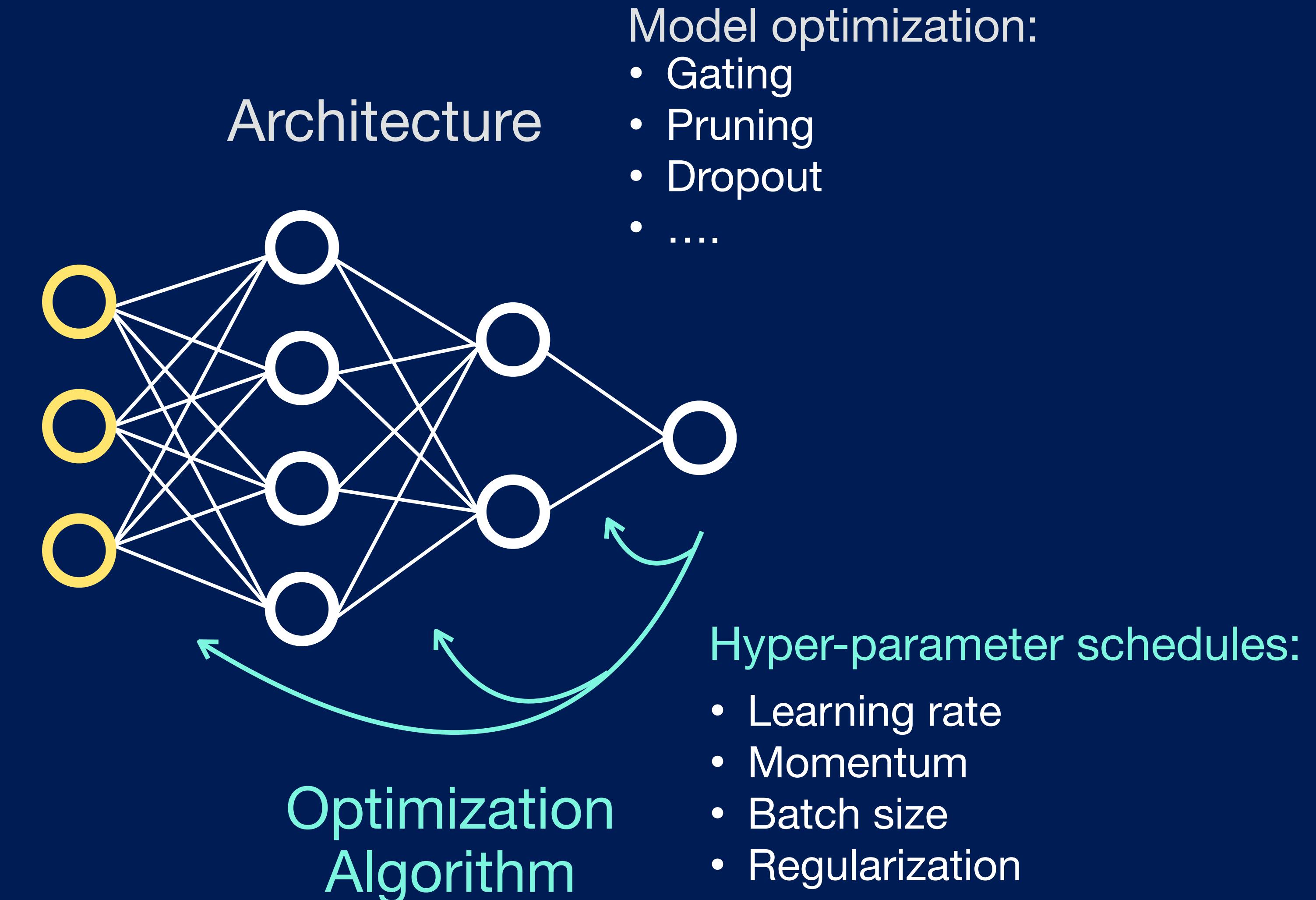


# Learning protocols for neural networks

Structured Data / Task

Dynamic data / task selection:

- Active learning
- Curriculum learning
- Transfer learning
- Multi-task learning
- ...



# Learning protocols for neural networks

**Complex interplay between time-dependent protocols and nonlinear learning dynamics**

Structured Data / Task

Dynamic data / task selection:

- Active learning
- Curriculum learning
- Transfer learning
- Multi-task learning
- ...



Model optimization:

- Gating
- Pruning
- Dropout
- ....

Hyper-parameter schedules:

- Learning rate
- Momentum
- Batch size
- Regularization
- ...

# \* Optimal learning protocols?

\* in the final performance

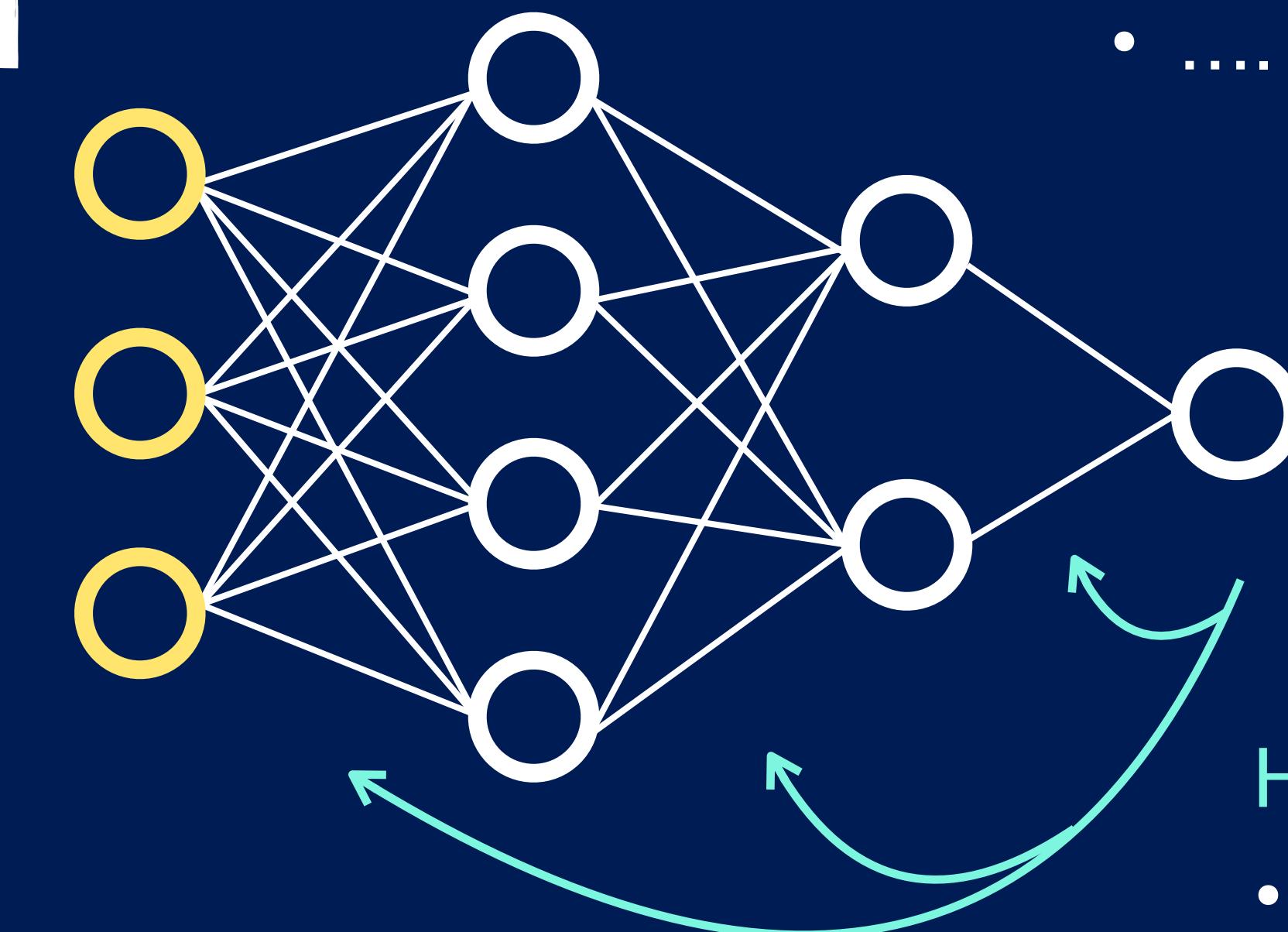
Complex interplay between time-dependent protocols and nonlinear learning dynamics

Structured Data / Task

Dynamic data / task selection:

- Active learning
- Curriculum learning
- Transfer learning
- Multi-task learning
- ...

Architecture



Optimization Algorithm

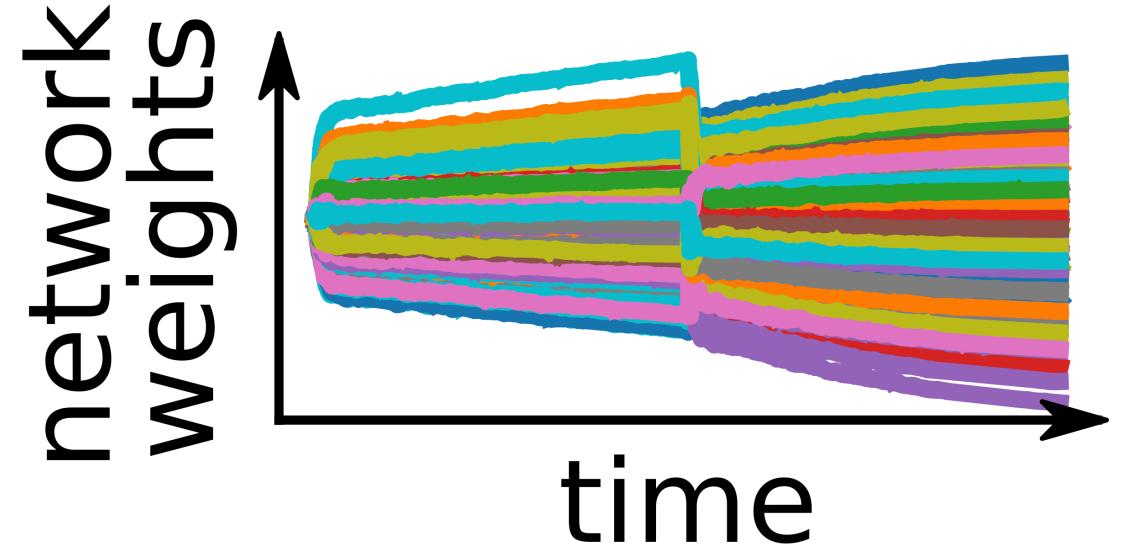
Model optimization:

- Gating
- Pruning
- Dropout
- ....

Hyper-parameter schedules:

- Learning rate
- Momentum
- Batch size
- Regularization
- ....

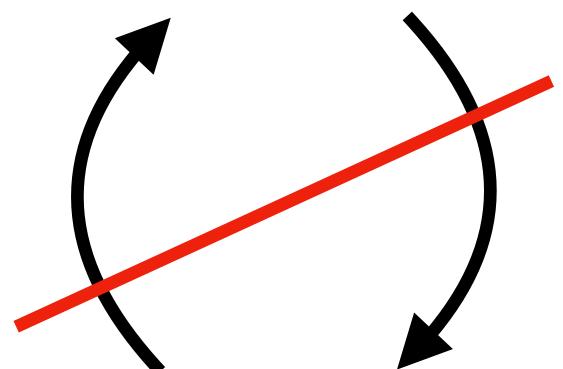
# Dimensionality reduction



$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu - \nabla \mathcal{L}^\mu(\mathbf{u}^\mu)$$

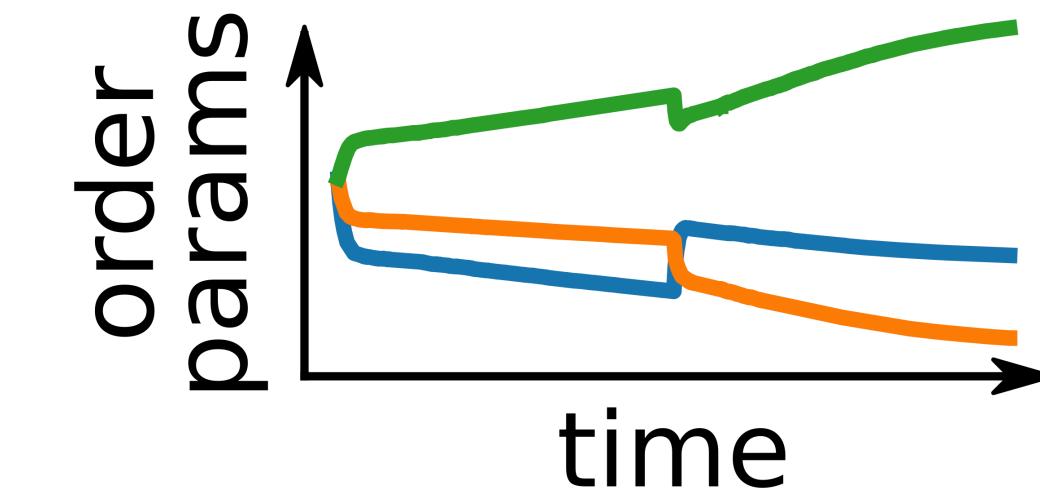
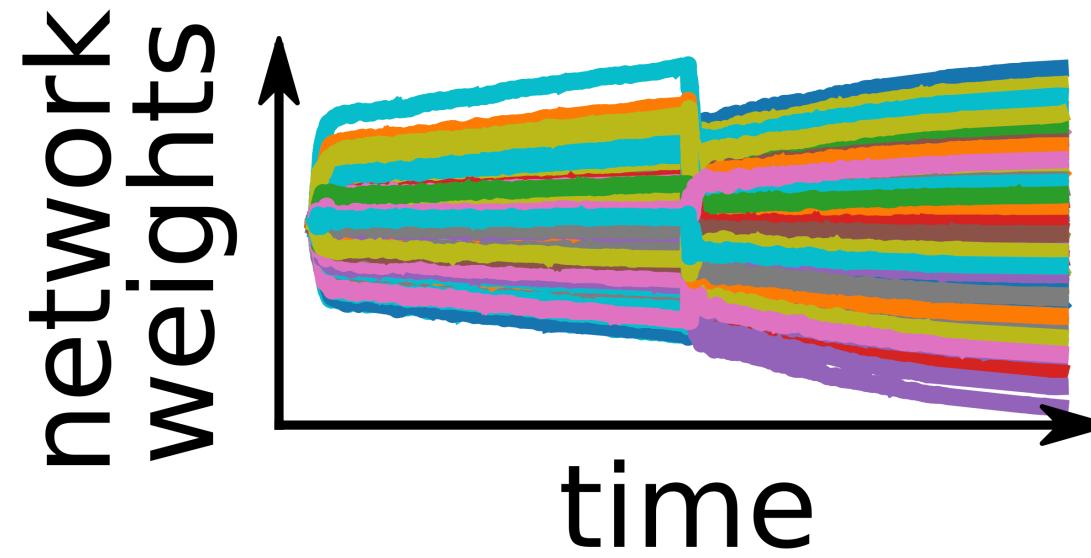
$\mathbf{w} \in \mathbb{R}^N$       control variables

High-dimensional learning dynamics



Control theory is computationally demanding

# Dimensionality reduction



$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu - \nabla \mathcal{L}^\mu(\mathbf{u}^\mu)$$

control variables

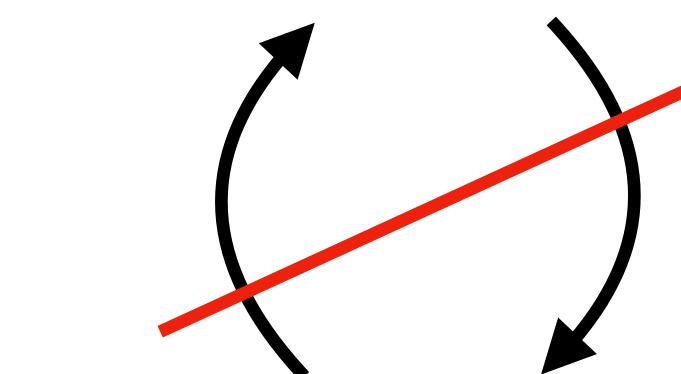
$$\mathbf{w} \in \mathbb{R}^N$$

High-dimensional learning dynamics

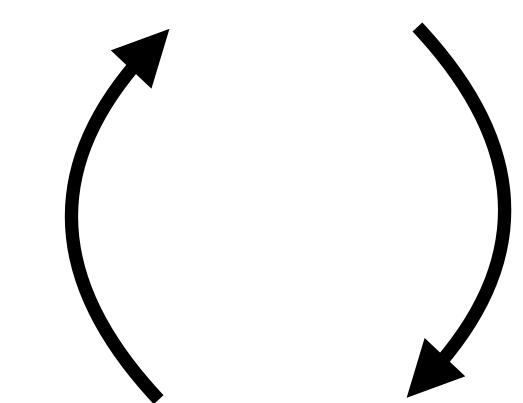
Statistical Physics

$$N \rightarrow \infty$$

effective dynamics of order parameters

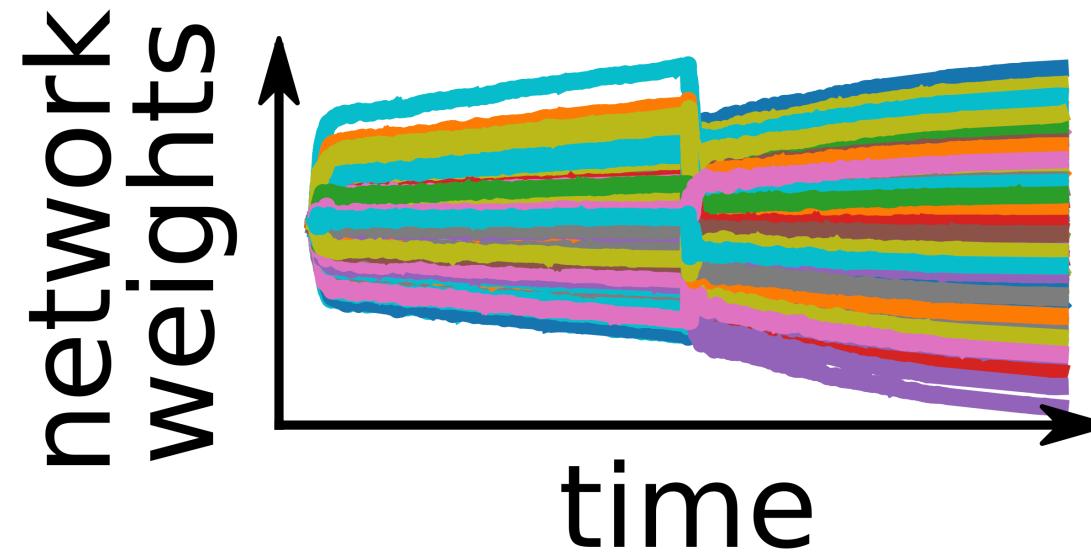


Control theory is computationally demanding



Control theory can be applied + interpretability

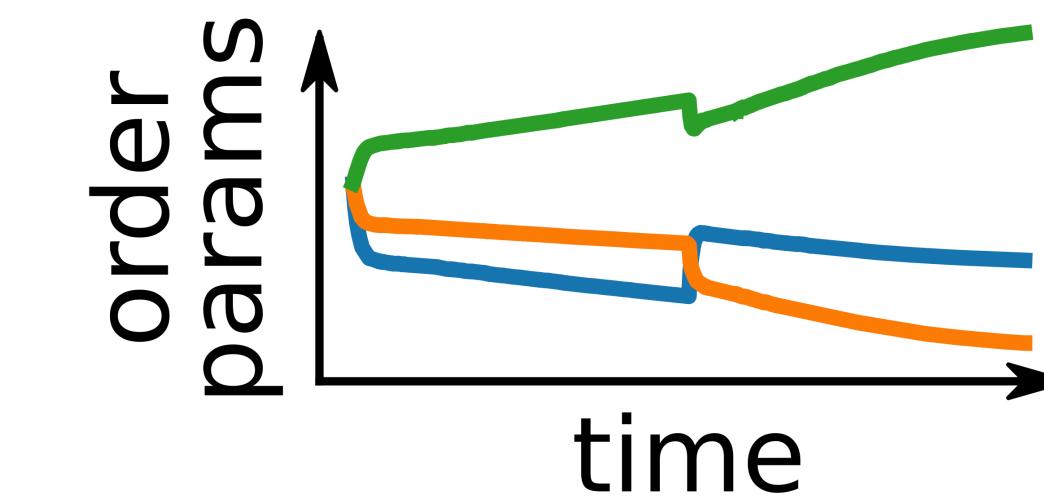
# Dimensionality reduction



$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu - \nabla \mathcal{L}^\mu(\mathbf{u}^\mu)$$

control variables

$$\mathbf{w} \in \mathbb{R}^N$$



$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

**Optimal control**

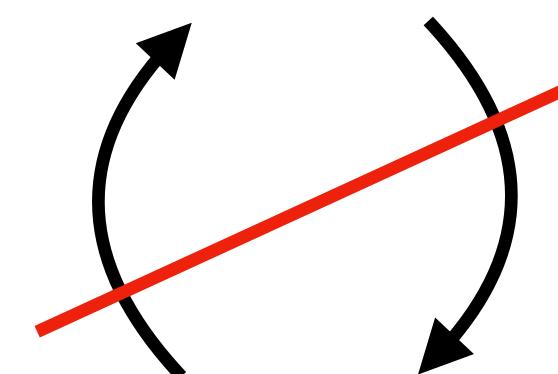
High-dimensional learning dynamics

Statistical Physics

$$N \rightarrow \infty$$

effective dynamics of order parameters

$$\mathbf{u}^*$$



Control theory is computationally demanding

Control theory can be applied + interpretability



Francesca Mignacco, FM,  
*arXiv:2507.07907*

# Setup – Data model

## Dataset

$$\mathcal{D} = \{\mathbf{x}_\mu, y_\mu\}_{\mu=1}^P$$

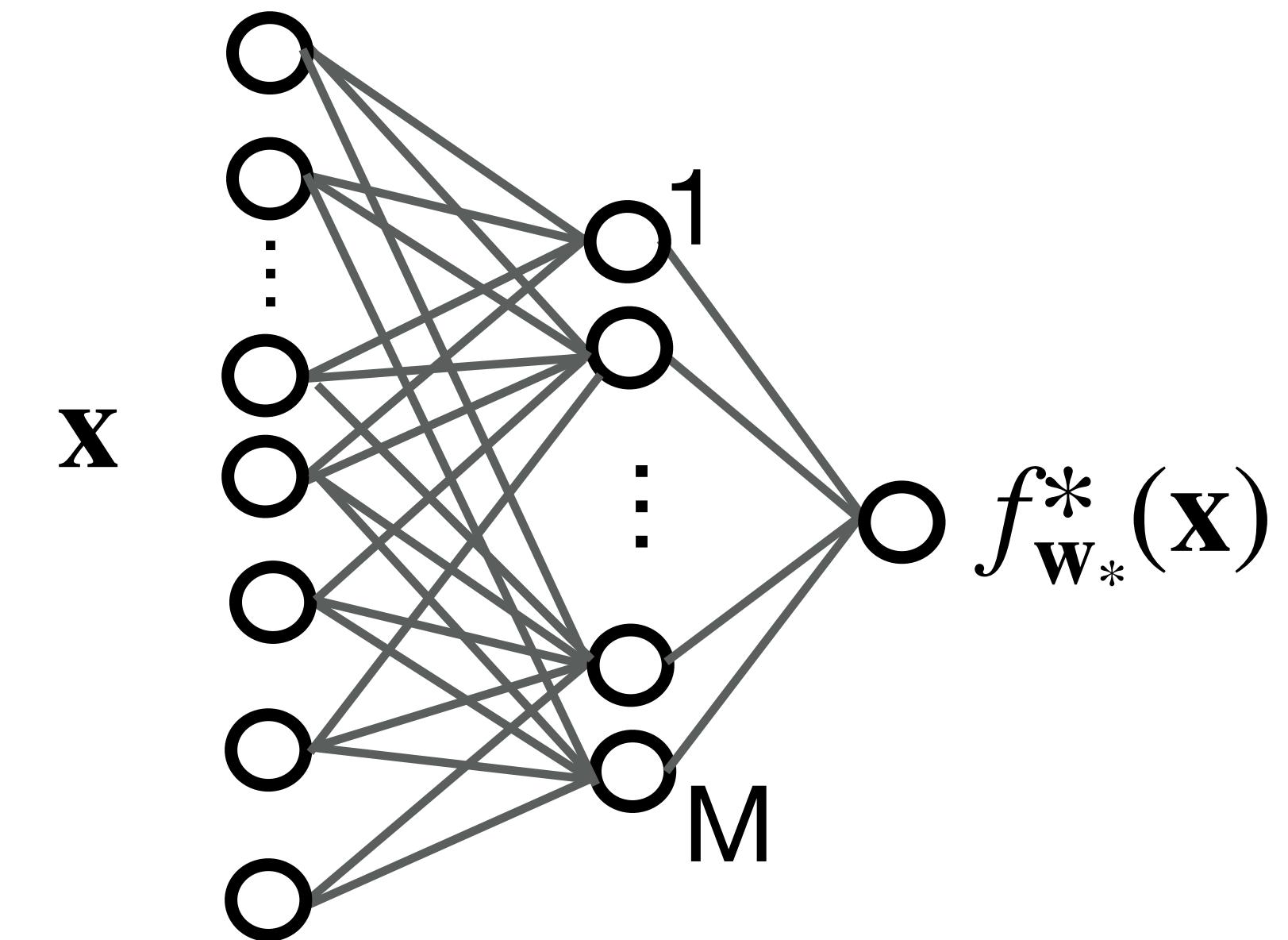
$$\mathbf{x}_\mu \sim \sum_{c=1}^C p_c \mathcal{N}\left(\frac{\mu_c}{\sqrt{N}}, \sigma_c^2 \mathbf{I}_N\right) \in \mathbb{R}^N$$

$$y_\mu = f_{\mathbf{w}_*}^*(\mathbf{x}_\mu) + \sigma_n z_\mu$$

Teacher Network  
 ( $M$  hidden nodes)

$$\mathbf{w}_* \in \mathbb{R}^{N \times M}$$

Label noise  
 $z_\mu \sim \mathcal{N}(0, 1)$

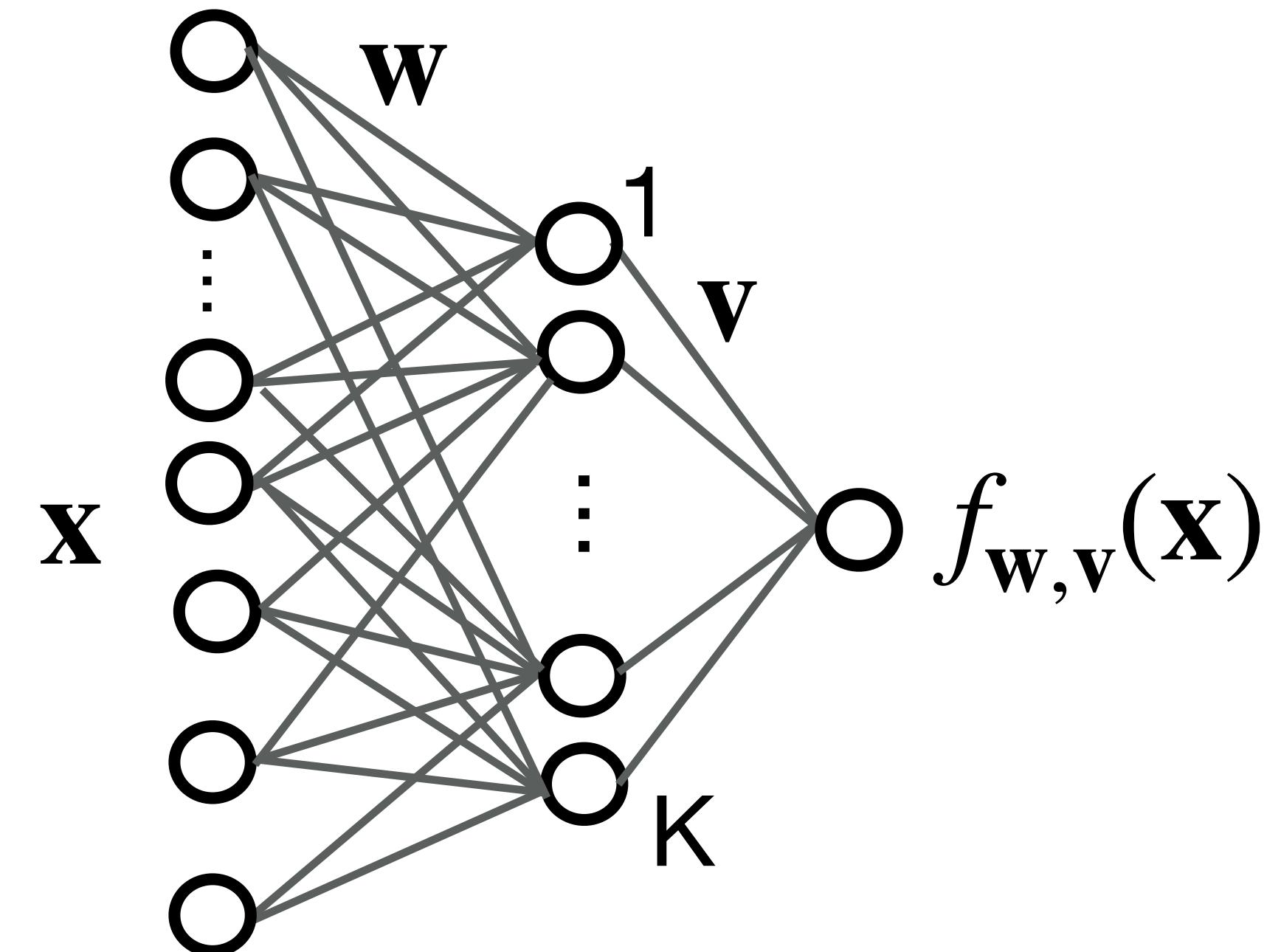


# Setup – Architecture and algorithm

Two-layer neural network  
 $(K$  hidden nodes)

$$f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) = \frac{1}{\sqrt{K}} \sum_{k=1}^K v_k g\left(\frac{\mathbf{w}_k \cdot \mathbf{x}}{\sqrt{N}}\right)$$

$$\mathbf{w} \in \mathbb{R}^{K \times M}$$



# Setup – Architecture and algorithm

**Two-layer neural network**  
 $(K$  hidden nodes)

$$f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) = \frac{1}{\sqrt{K}} \sum_{k=1}^K v_k g\left(\frac{\mathbf{w}_k \cdot \mathbf{x}}{\sqrt{N}}\right)$$

$$\mathbf{w} \in \mathbb{R}^{K \times M}$$

**Per-sample cost function**

$$\mathcal{L}(\mathbf{w}, \mathbf{v} \mid \mathbf{x}, y) = \ell(f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}), y) + g(\mathbf{w}, \mathbf{v})$$

Loss function

Regularization  
function

# Setup – Architecture and algorithm

**Two-layer neural network**  
 $(K$  hidden nodes)

$$f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) = \frac{1}{\sqrt{K}} \sum_{k=1}^K v_k g\left(\frac{\mathbf{w}_k \cdot \mathbf{x}}{\sqrt{N}}\right)$$

$$\mathbf{w} \in \mathbb{R}^{K \times M}$$

**Per-sample cost function**

$$\mathcal{L}(\mathbf{w}, \mathbf{v} \mid \mathbf{x}, y) = \ell(f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}), y) + g(\mathbf{w}, \mathbf{v})$$

Loss function

Regularization  
function

**One-pass SGD**

$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu - \eta \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^\mu, \mathbf{v}^\mu \mid \mathbf{x}^\mu, y^\mu)$$

# Dynamics of order parameters

$$\mathbf{S} \in \mathbb{R}^{K \times K}$$

$$\mathbf{M} \in \mathbb{R}^{K \times M}$$

$$\mathbf{R} \in \mathbb{R}^{K \times C}$$

$$S_{kk'} = \frac{\mathbf{w}_k \cdot \mathbf{w}_{k'}}{N}$$

$$M_{km} = \frac{\mathbf{w}_k \cdot \mathbf{w}_{*,m}}{N}$$

$$R_{kc} = \frac{\mathbf{w}_k \cdot \boldsymbol{\mu}_c}{N}$$

$$\mathbf{Q} \equiv (\text{vec}(\mathbf{S}), \text{vec}(\mathbf{M}), \text{vec}(\mathbf{R})) \in \mathbb{R}^{K(K+M+C)}$$

One-pass SGD

$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu - \eta \nabla_{\mathbf{w}} \mathcal{L}^\mu(\mathbf{u}_\mu) \xrightarrow[N \rightarrow \infty, \mu \rightarrow \infty]{\text{with } \alpha = \mu/N} \frac{d\mathbf{Q}(\alpha)}{d\alpha}$$

$$\mu = 1, \dots, P$$

ODEs for the order parameters

$$\frac{d\mathbf{Q}(\alpha)}{d\alpha} = f_{\mathbf{Q}}(\mathbf{Q}(\alpha), \mathbf{u}(\alpha))$$

$$\alpha \in [0, \alpha_F]$$

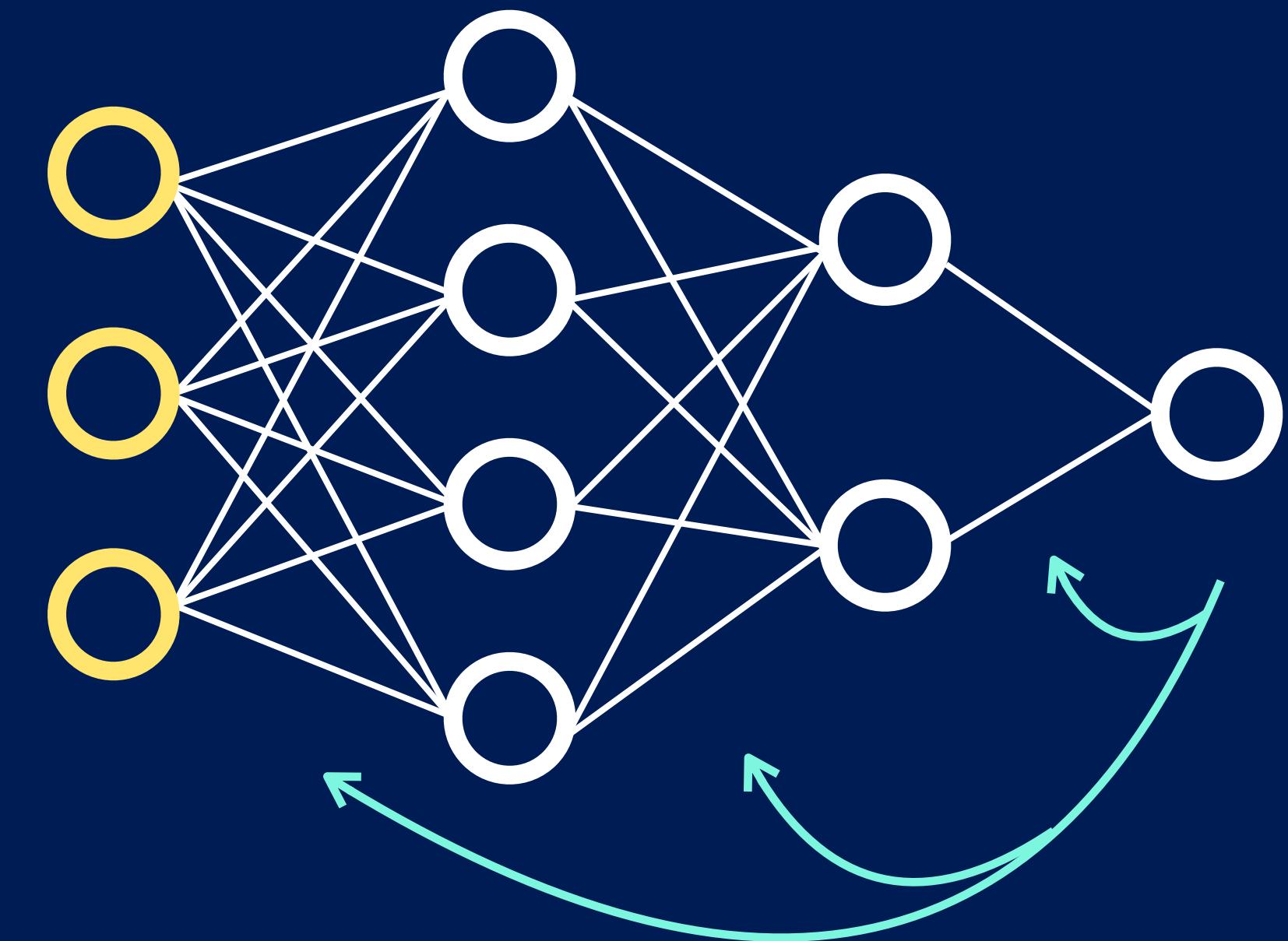
[Saad & Solla (1995); Biehl & Schwarze (1995); Goldt et al (2019); Veiga et al. (2022), Arnaboldi et al (2023) ... and many more]

# Optimal learning protocols



Francesca  
Mignacco  
CUNY & Princeton

Structured  
Data / Task



Optimization  
Algorithm

**FM**, F. Mignacco, *arXiv:2505.07792*

F. Mignacco, **FM**, *arXiv:2507.07907*



# Optimal learning protocols



Francesca  
Mignacco  
CUNY & Princeton



Stefano Sarao Mannelli  
Chalmers University  
Gothenburg University

Structured  
Data / Task

2

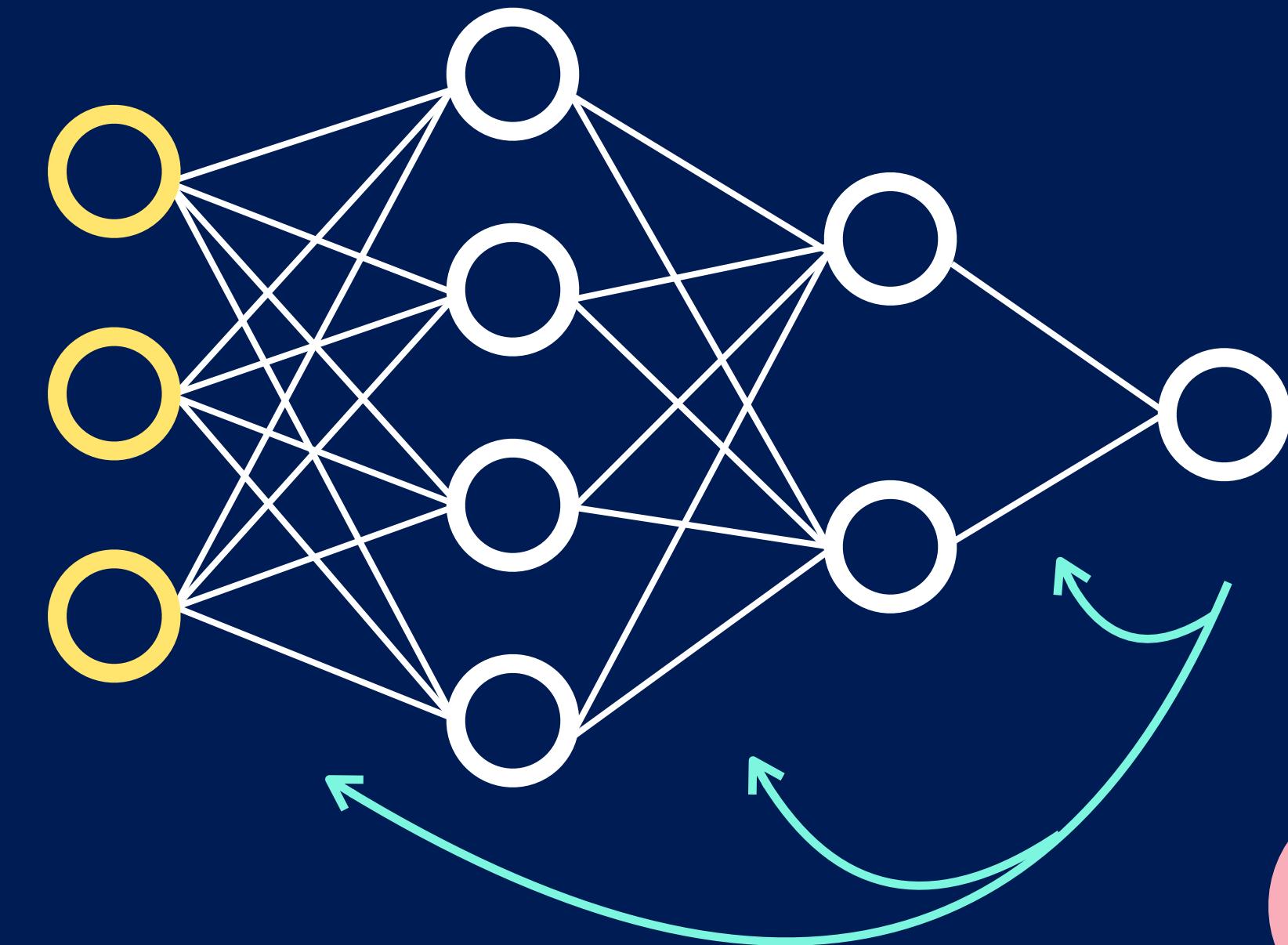
Continual learning

FM, F. Mignacco, *arXiv:2505.07792*

F. Mignacco, FM, *arXiv:2507.07907*

FM, S. Sarao Mannelli, F. Mignacco, *ICLR 2025*

Architecture



Optimization  
Algorithm

1

Dropout

2

Learning rate

# Optimal learning protocols



Francesca  
Mignacco  
CUNY & Princeton



Stefano Sarao Mannelli  
Chalmers University  
Gothenburg University

Structured  
Data / Task

1

Continual learning

2

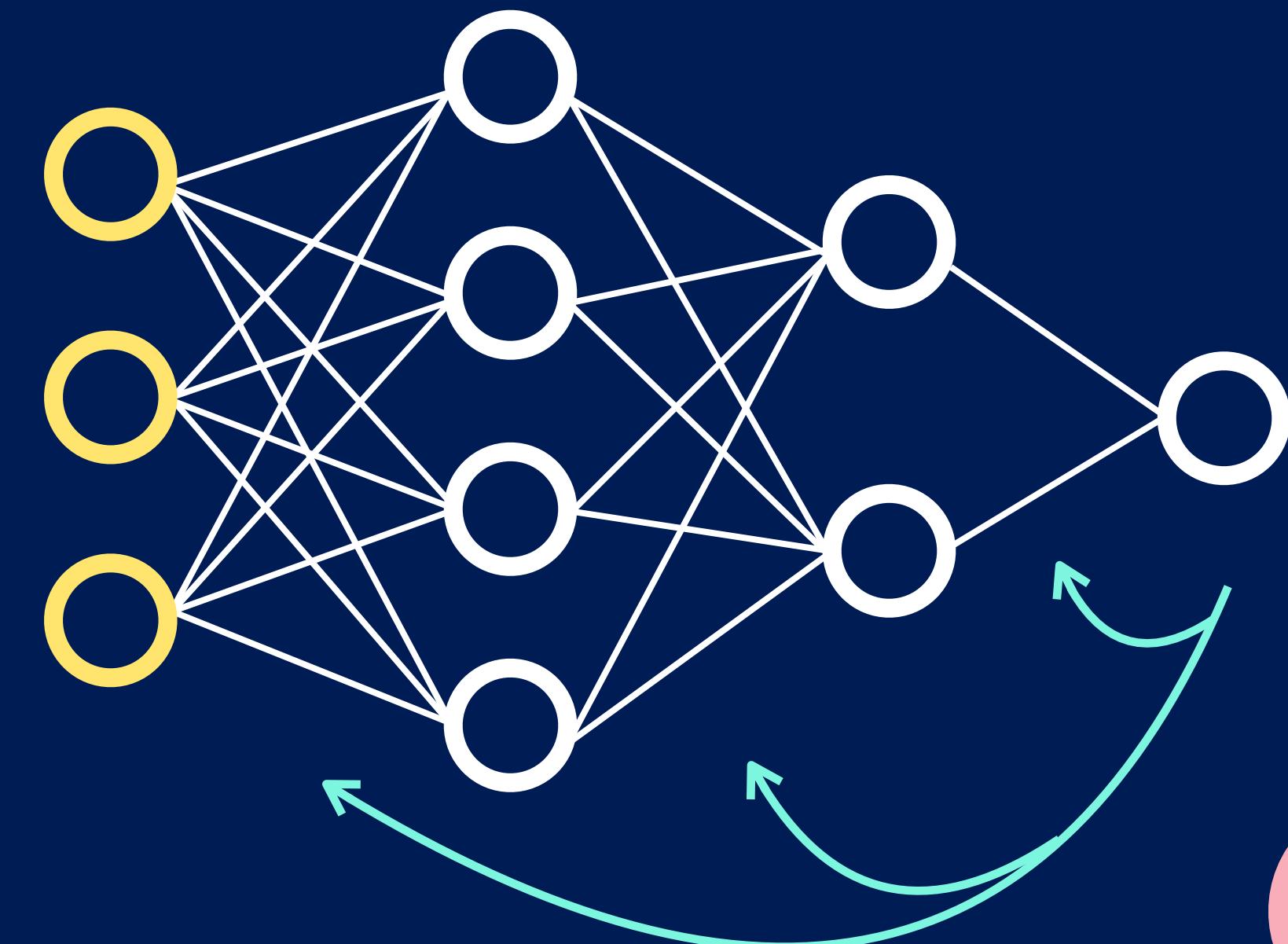
Denoising autoencoder

FM, F. Mignacco, *arXiv:2505.07792*

F. Mignacco, FM, *arXiv:2507.07907*

FM, S. Sarao Mannelli, F. Mignacco, *ICLR 2025*

Architecture



Optimization  
Algorithm

1

Dropout

2

Learning rate

3

Batch size

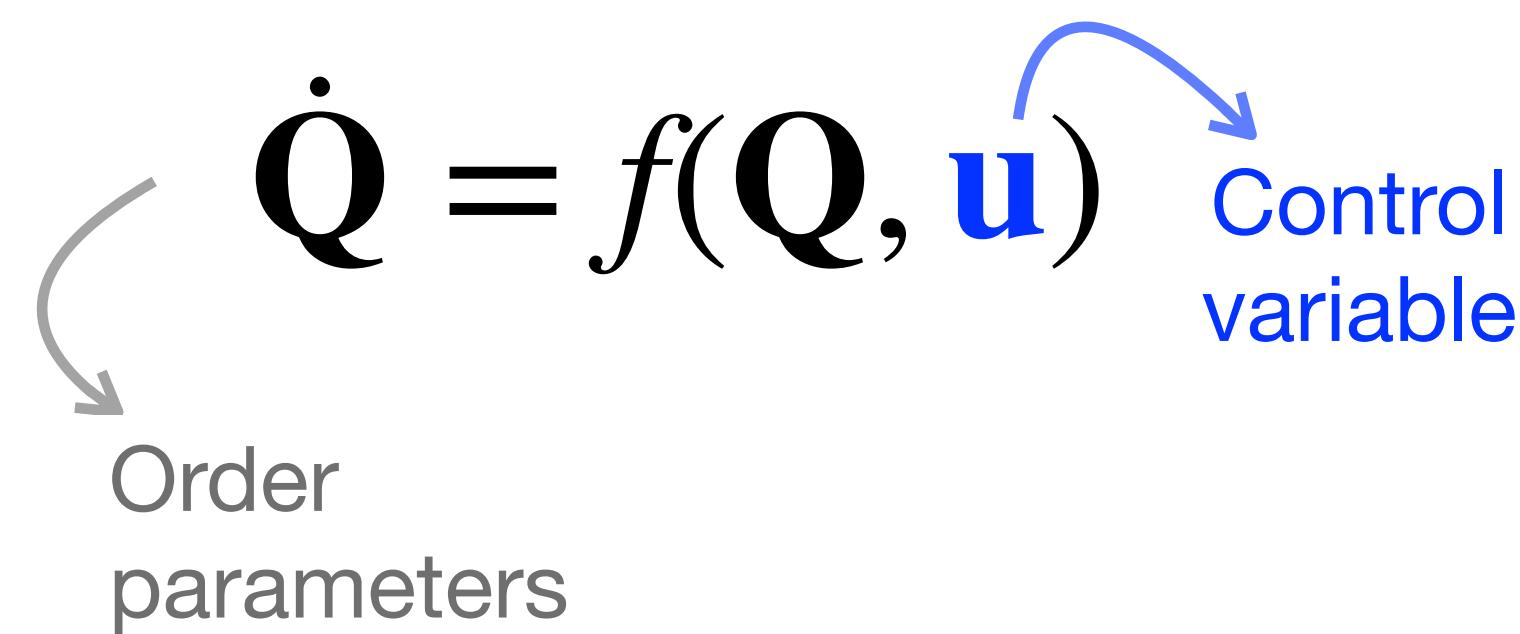
# Can we compute optimal meta-parameter schedules?

Forward learning dynamics

$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

Order parameters

Control variable



# Pontryagin's maximum principle

Variational approach

Cost functional:

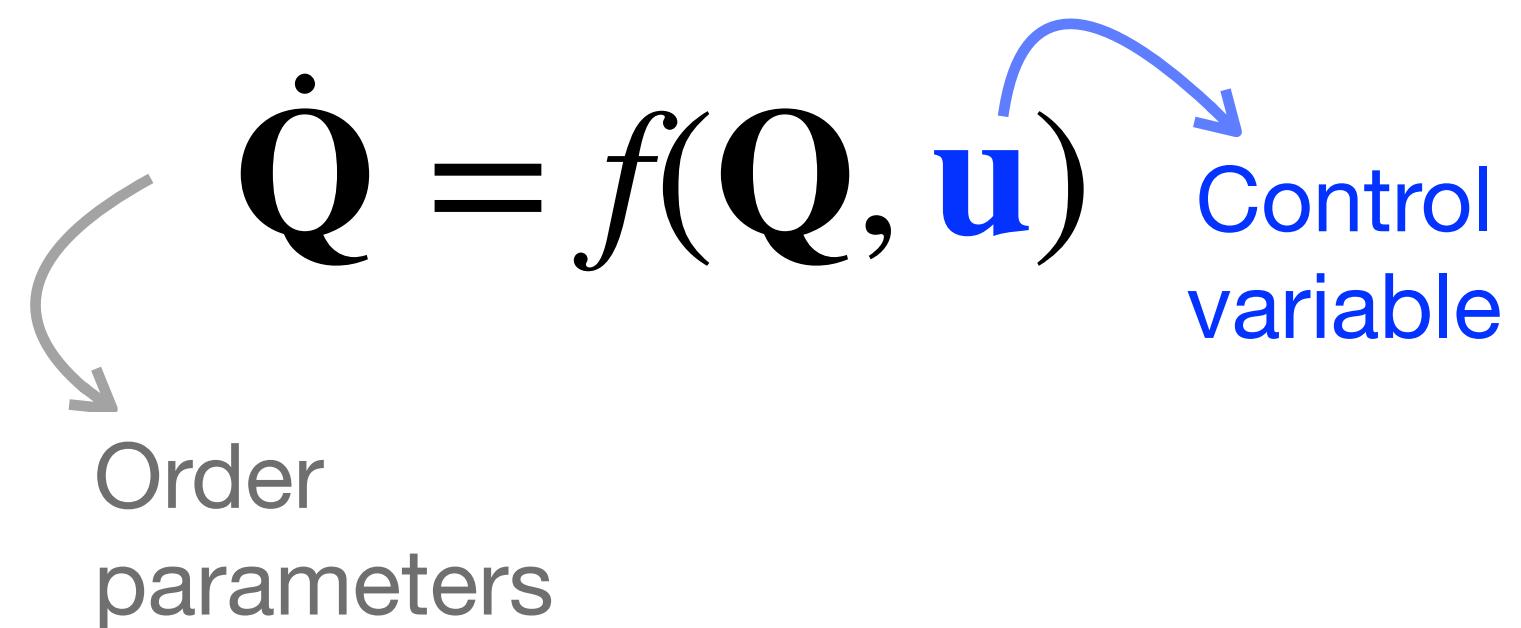
$$\mathcal{F}[\mathbf{Q}, \mathbf{u}] = \langle \epsilon_g(\alpha_F) \rangle$$

## Forward learning dynamics

$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

Order parameters

Control variable



# Pontryagin's maximum principle

Variational approach  
 Cost functional:

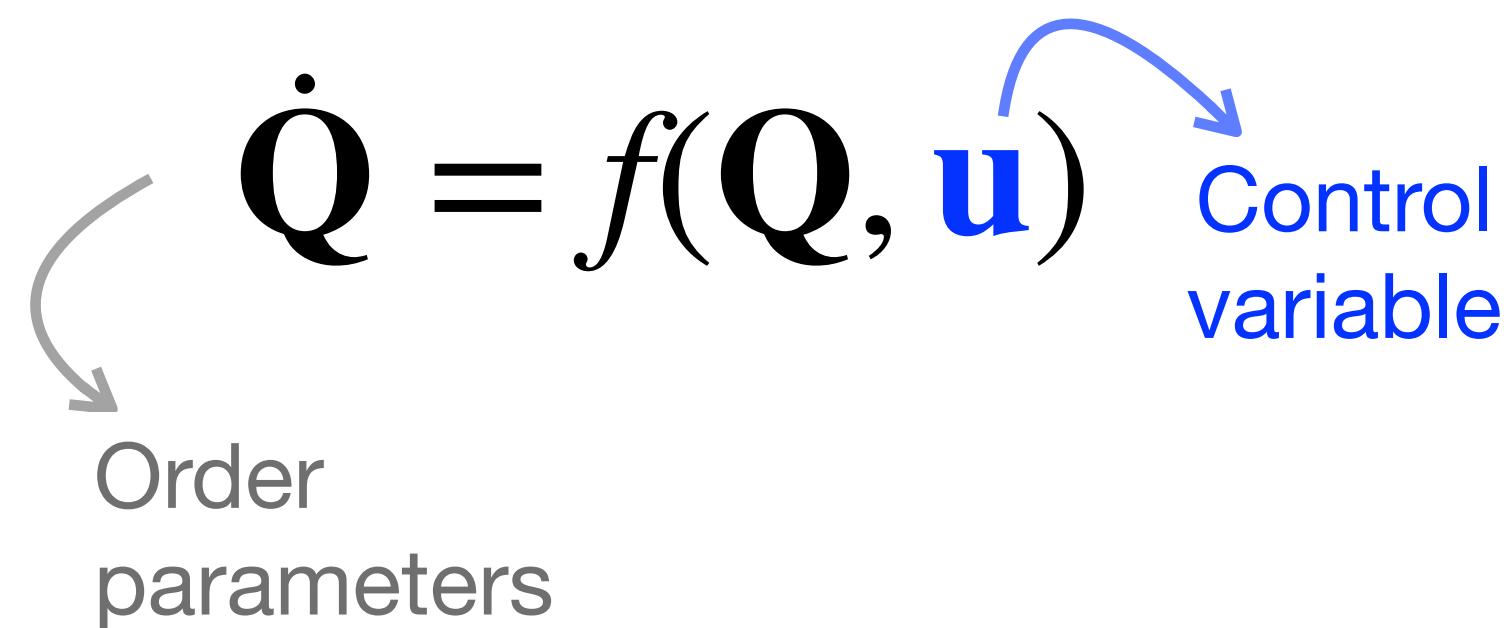
$$\mathcal{F}[\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{u}] = \langle \epsilon_g(\alpha_F) \rangle + \int_0^{\alpha_F} d\alpha \hat{\mathbf{Q}}(\alpha) \cdot [-\dot{\mathbf{Q}}(\alpha) + f(\mathbf{Q}(\alpha), \mathbf{u}(\alpha))]$$

Forward learning dynamics

$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

Order parameters

Control variable



# Pontryagin's maximum principle

Variational approach  
 Cost functional:

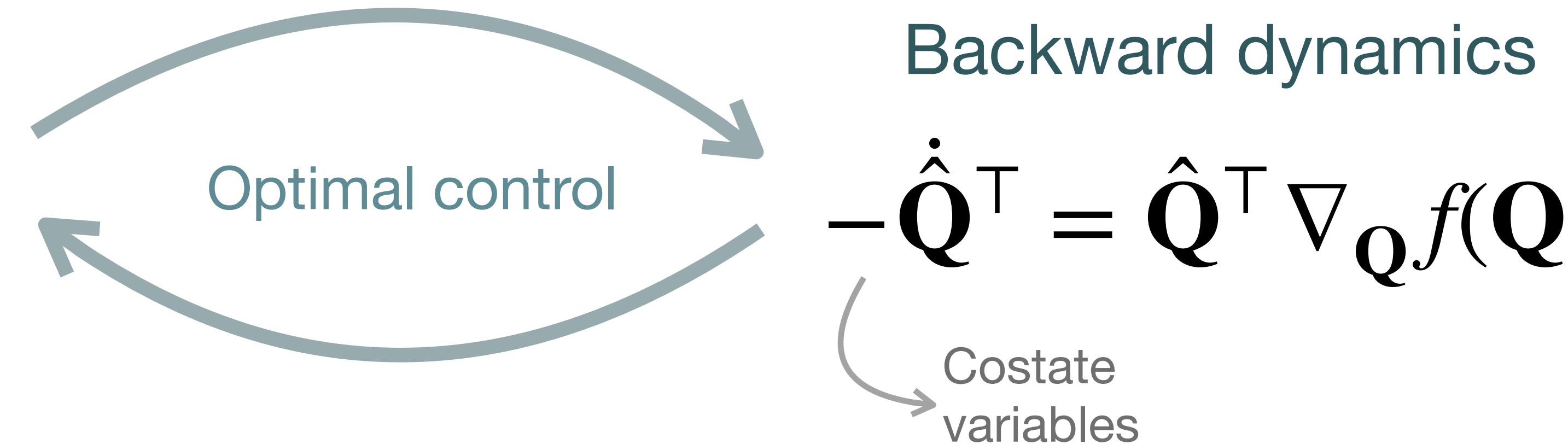
$$\mathcal{F}[\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{u}] = \langle \epsilon_g(\alpha_F) \rangle + \int_0^{\alpha_F} d\alpha \hat{\mathbf{Q}}(\alpha) \cdot [-\dot{\mathbf{Q}}(\alpha) + f(\mathbf{Q}(\alpha), \mathbf{u}(\alpha))]$$

Forward learning dynamics

$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

Control variable

Order parameters



$$-\dot{\hat{\mathbf{Q}}^\top} = \hat{\mathbf{Q}}^\top \nabla_{\mathbf{Q}} f(\mathbf{Q}, \mathbf{u})$$

# Pontryagin's maximum principle

Variational approach  
 Cost functional:

$$\mathcal{F}[\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{u}] = \langle \epsilon_g(\alpha_F) \rangle + \int_0^{\alpha_F} d\alpha \hat{\mathbf{Q}}(\alpha) \cdot [-\dot{\mathbf{Q}}(\alpha) + f(\mathbf{Q}(\alpha), \mathbf{u}(\alpha))]$$

Forward learning dynamics

$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u})$$

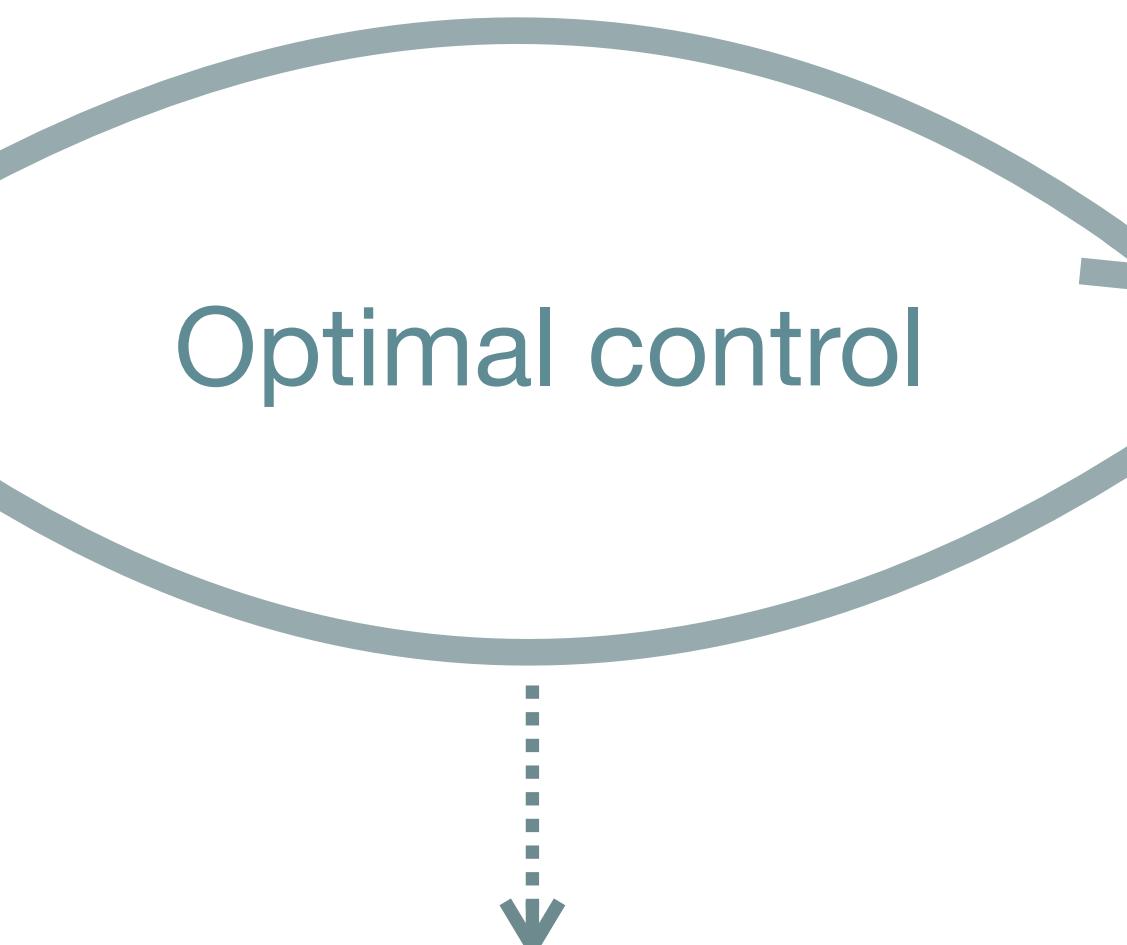
Control variable

Order parameters

Backward dynamics

$$-\dot{\hat{\mathbf{Q}}}^\top = \hat{\mathbf{Q}}^\top \nabla_{\mathbf{Q}} f(\mathbf{Q}, \mathbf{u})$$

Costate variables



Optimality condition

$$\mathbf{u}^*(\alpha) = \operatorname{argmin}_{\mathbf{u}(\alpha)} \hat{\mathbf{Q}}(\alpha)^\top f(\mathbf{Q}(\alpha), \mathbf{u}(\alpha))$$

# Part I: *Dropout Regularization*

**FM**, F. Mignacco, *arXiv:2505.07792*

F. Mignacco, **FM**, *arXiv:2507.07907*

# Dropout regularization

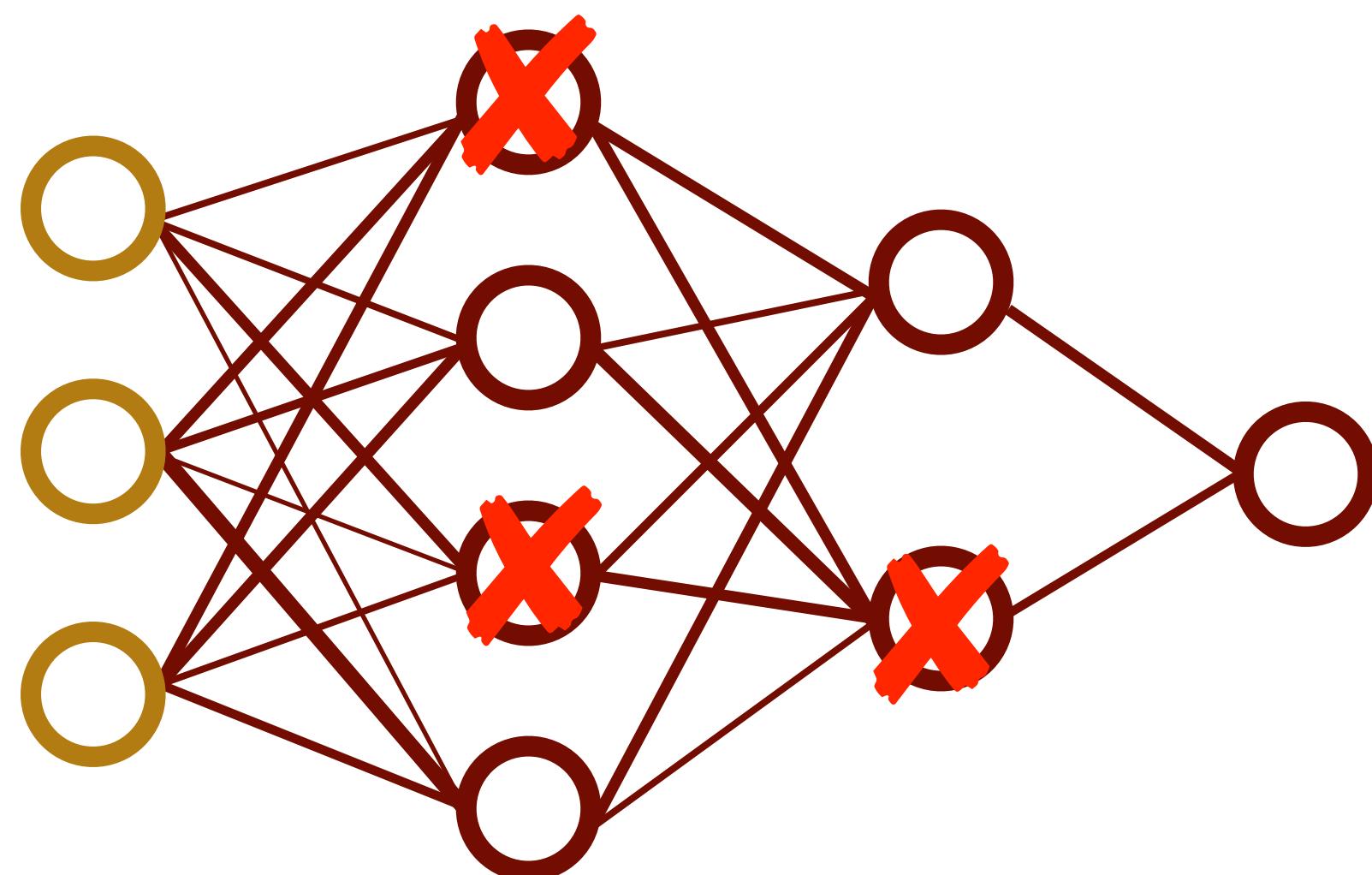
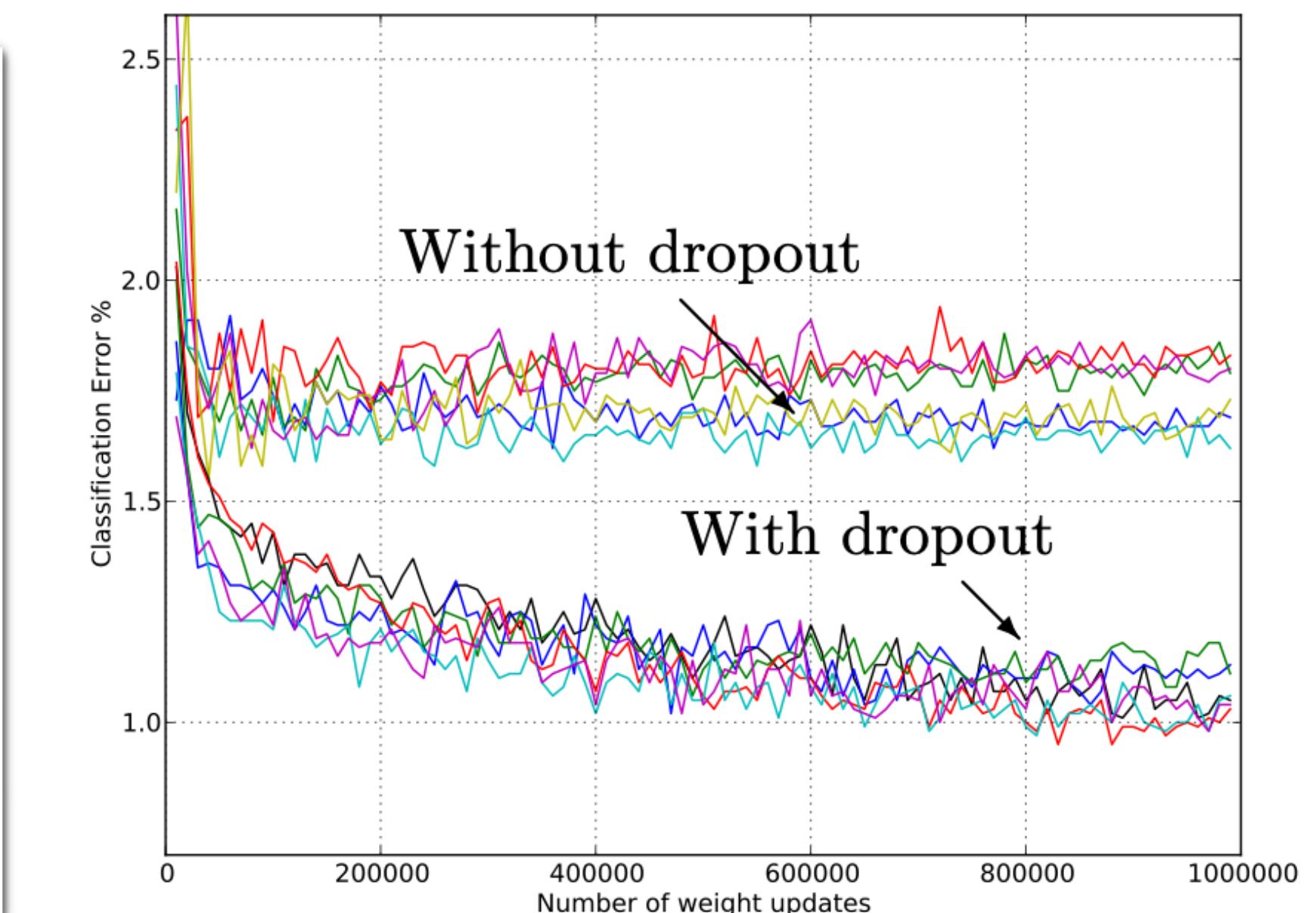
Journal of Machine Learning Research 15 (2014) 1929-1958

Submitted 11/13; Published 6/14

## Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Nitish Srivastava  
 Geoffrey Hinton  
 Alex Krizhevsky  
 Ilya Sutskever  
 Ruslan Salakhutdinov

NITISH@CS.TORONTO.EDU  
 HINTON@CS.TORONTO.EDU  
 KRIZ@CS.TORONTO.EDU  
 ILYA@CS.TORONTO.EDU  
 RSALAKHU@CS.TORONTO.EDU



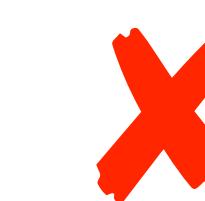
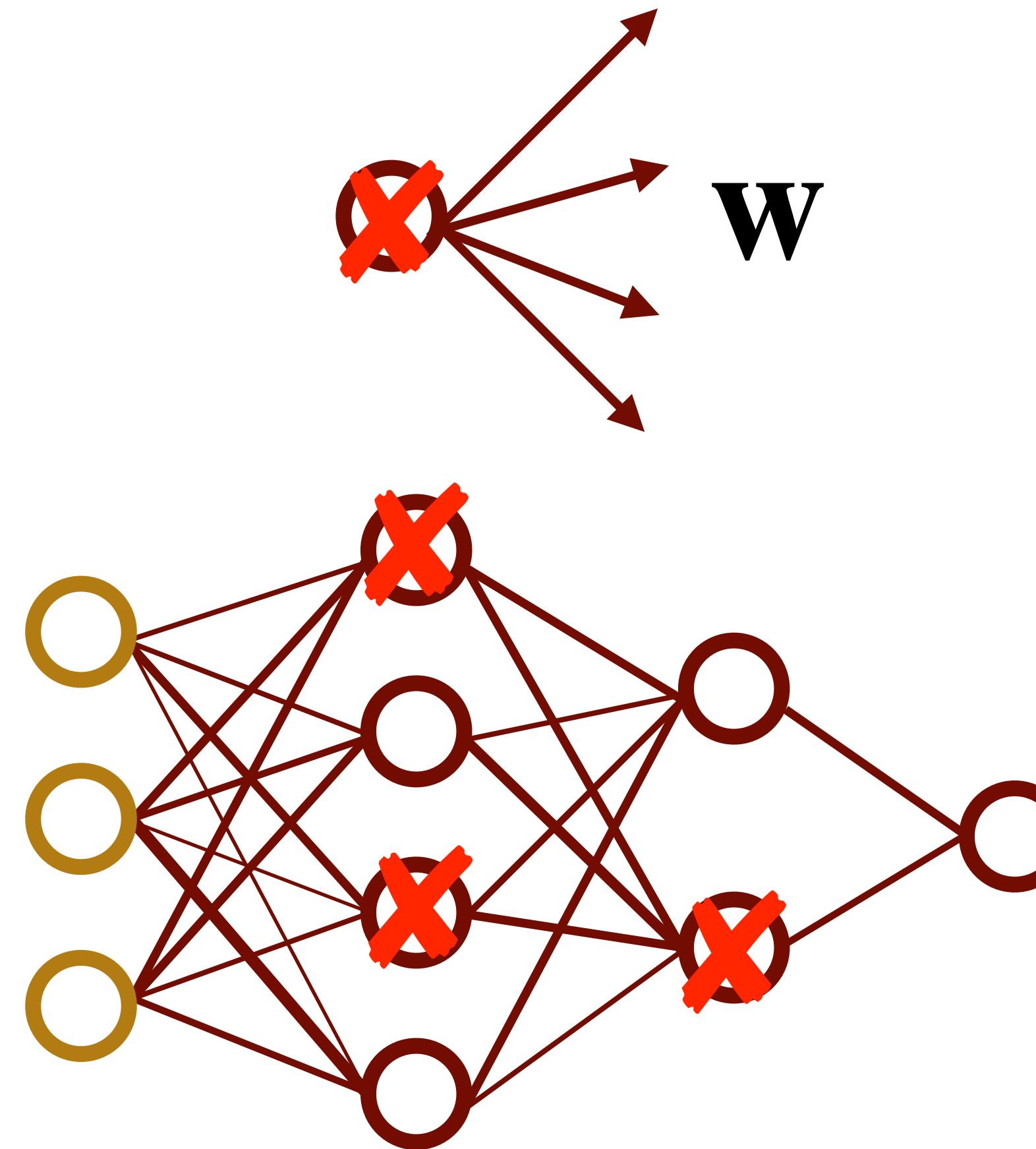
 = remove node with probability  $1 - p$

Figure 4: Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.

# Dropout regularization

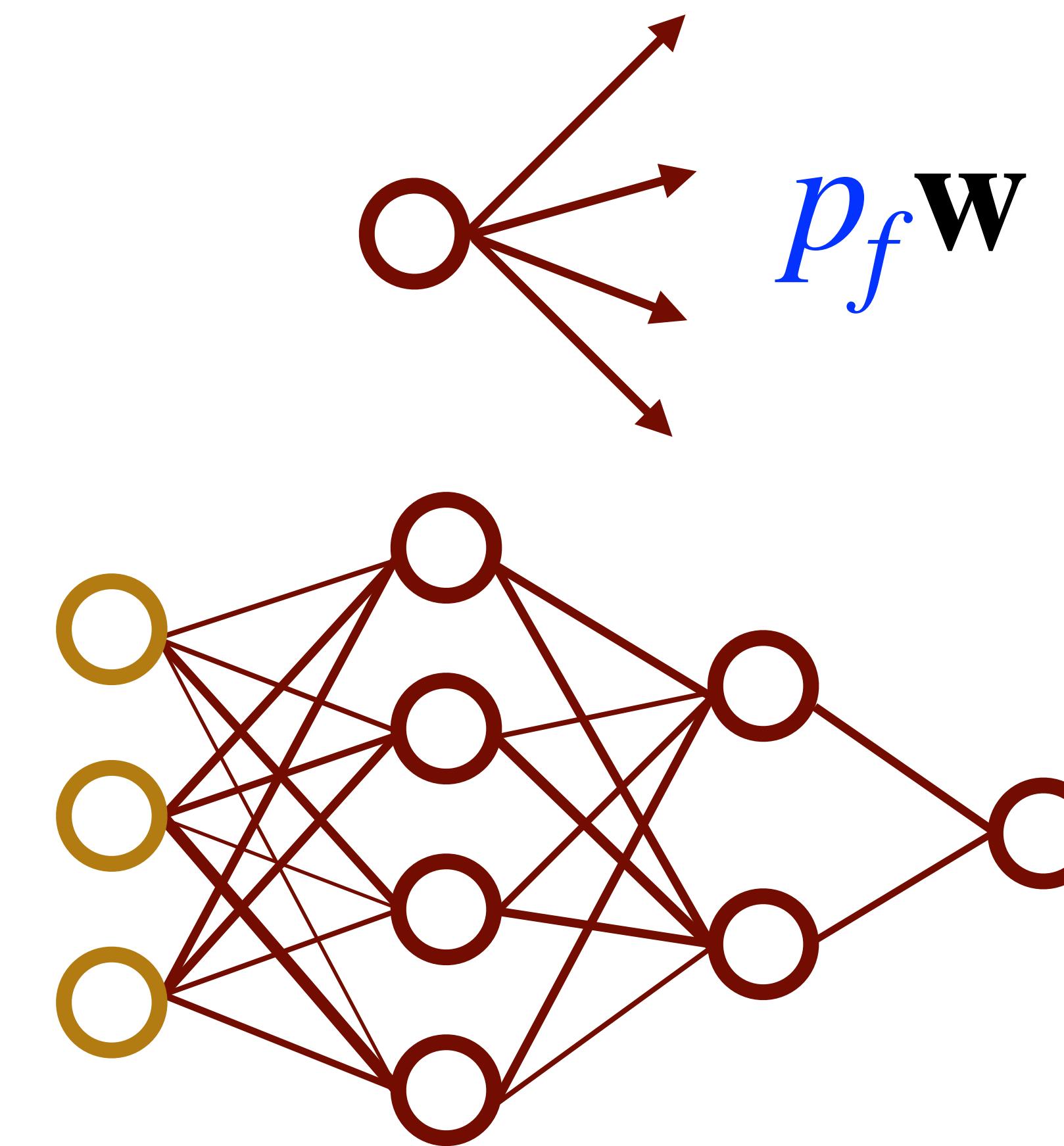
**Train**

✗ = remove node (with probability  $1 - p$ )



**Test**

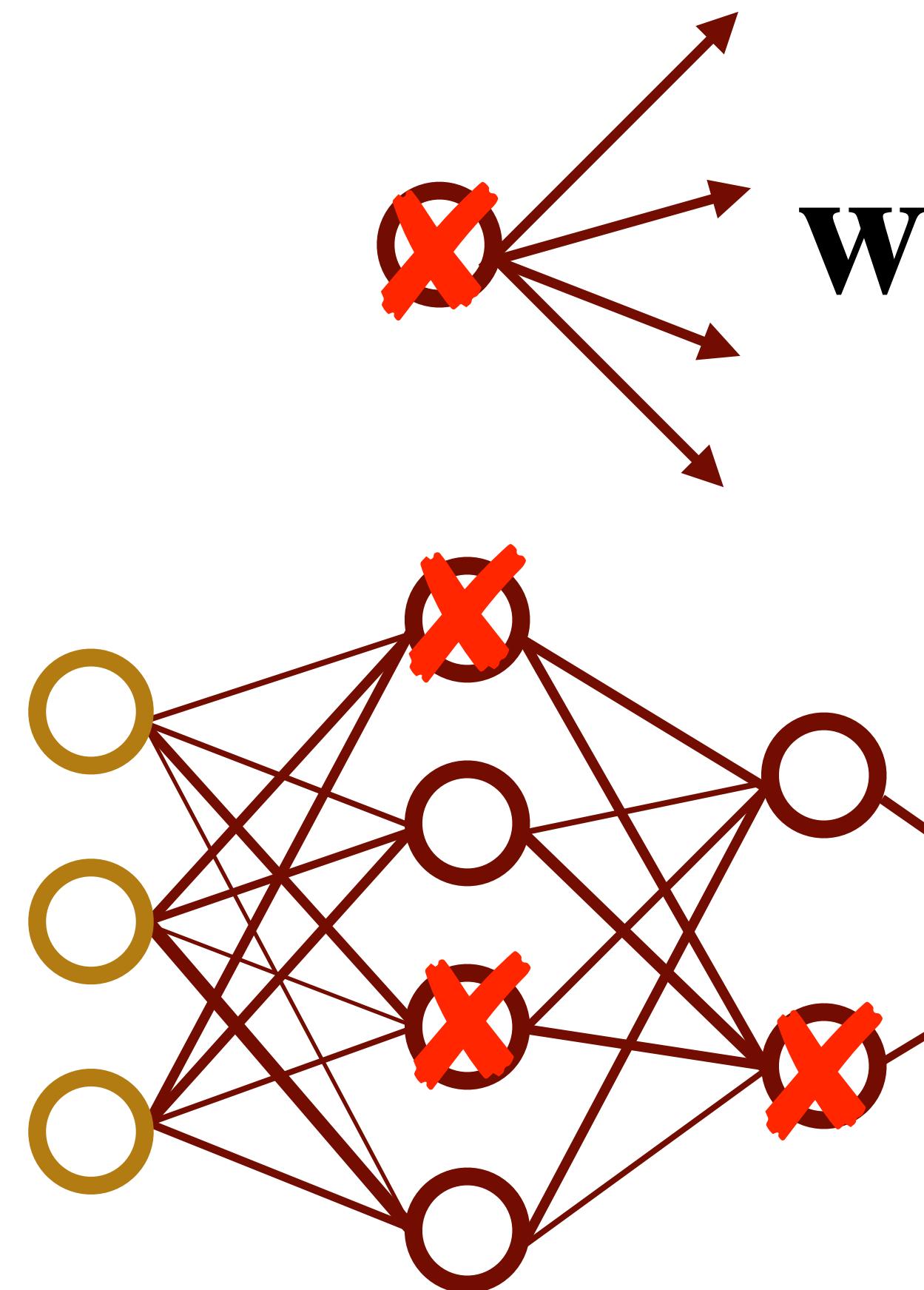
Always present but *rescaled* weights



# What is the *optimal* dropout probability?

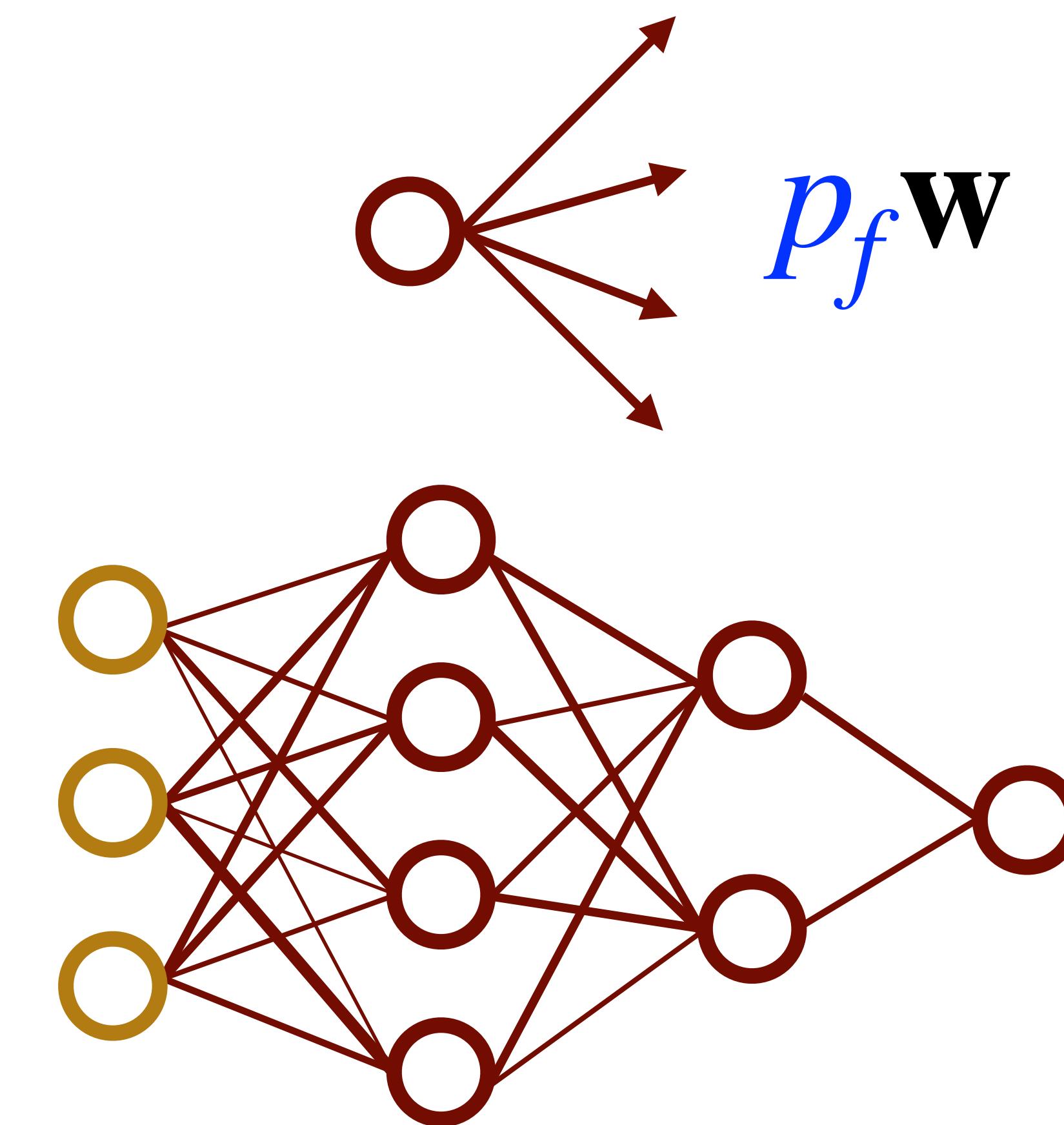
## Train

✗ = remove node (with probability  $1 - p$ )



## Test

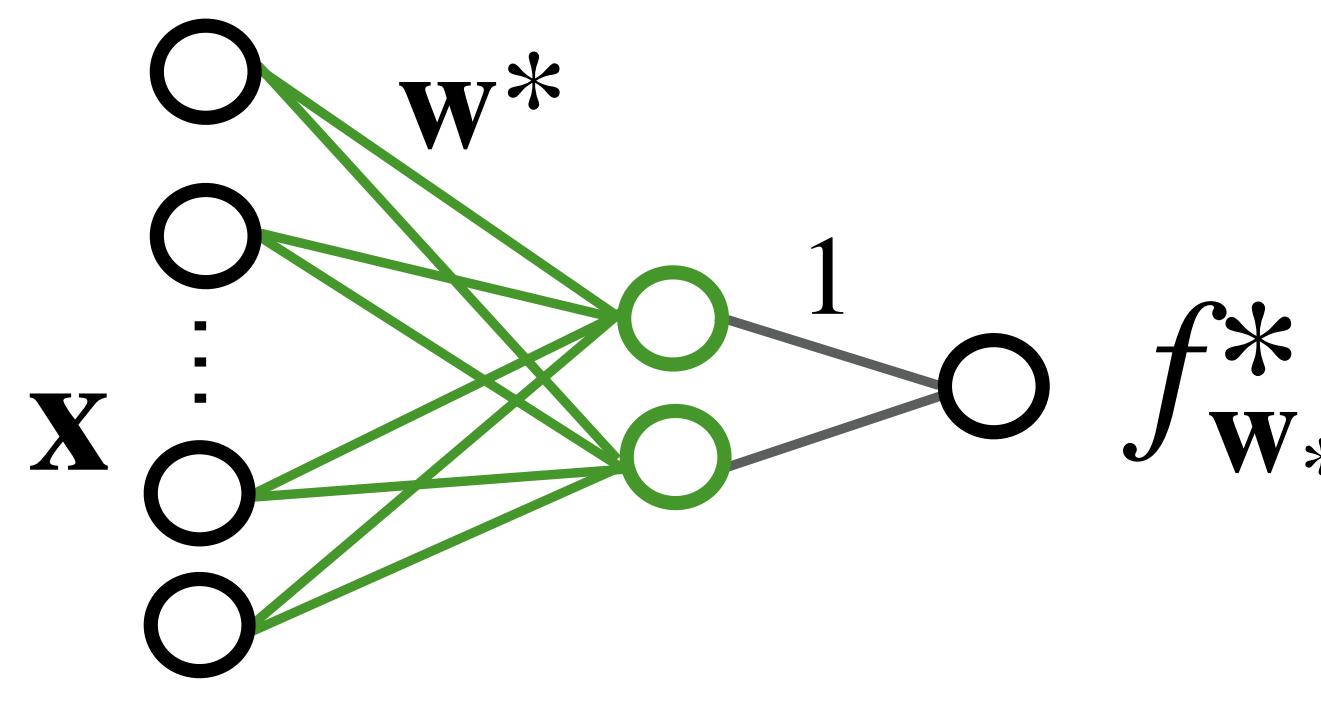
Always present but *rescaled* weights



# A teacher-student model of dropout

$$\mathbf{x} \sim \mathcal{N}(0, I_N)$$

**Teacher**

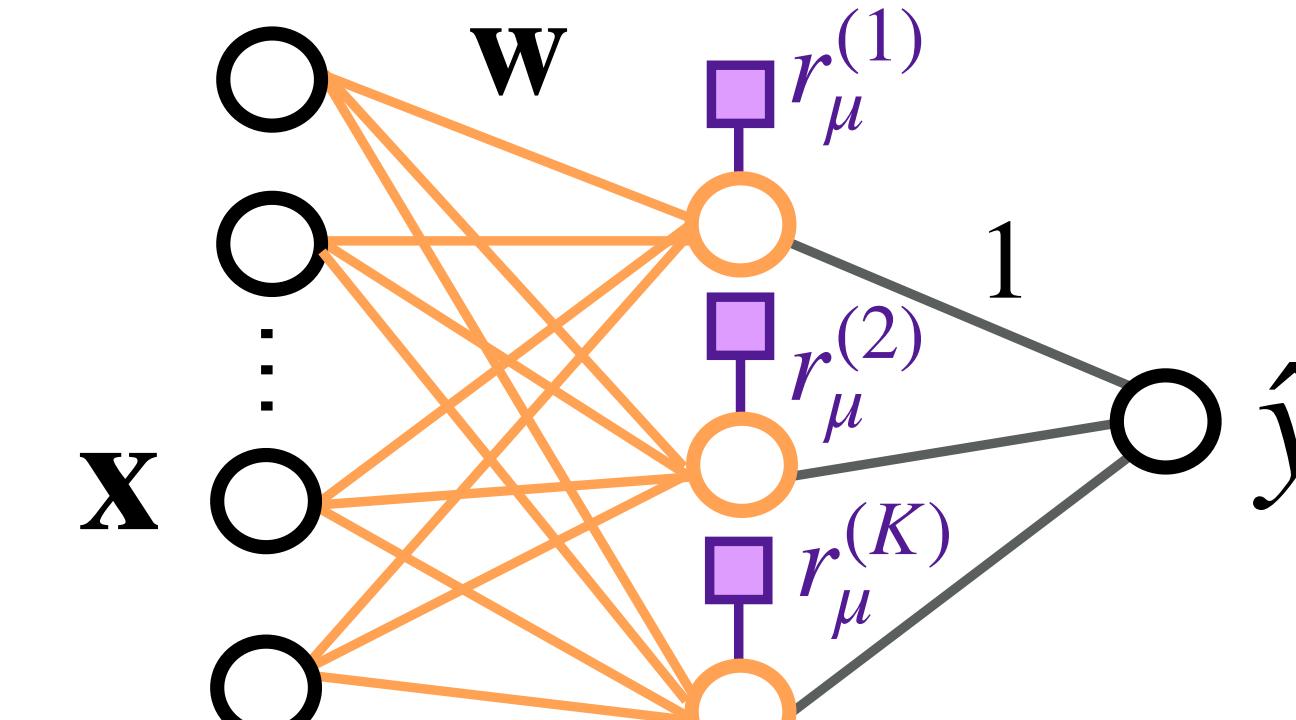


$$y = f_{\mathbf{w}^*}^*(\mathbf{x}) + \sigma_n z$$

Label noise:  $z \sim \mathcal{N}(0, 1)$

**Student**

(at training step  $\mu$ )



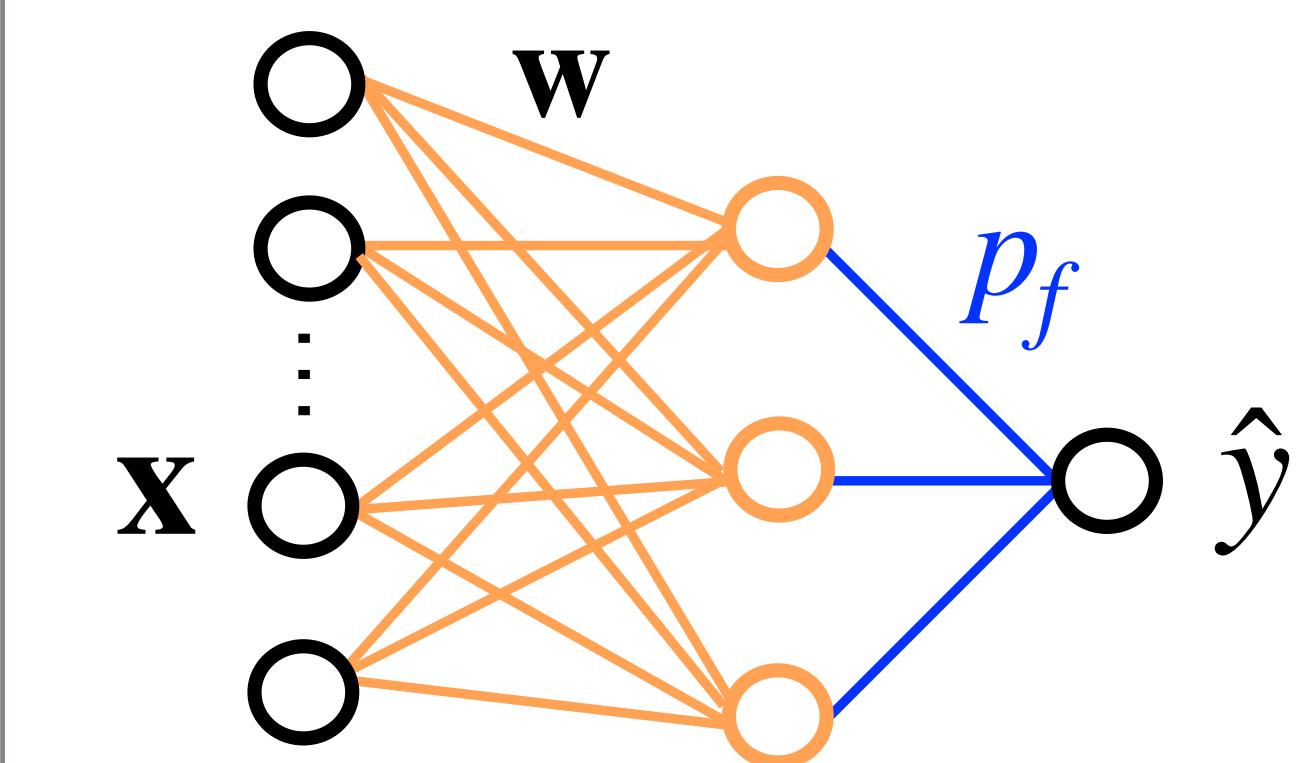
$K$  hidden nodes

*Node-activation variables*

$$r_\mu^{(1)}, r_\mu^{(2)}, \dots, r_\mu^{(K)} \sim \text{Bernoulli}(p_\mu)$$

**Student**

(at testing time)



$K$  hidden nodes

*Rescaling factor:  $p_f$*

# Online learning with constant node-activation prob.

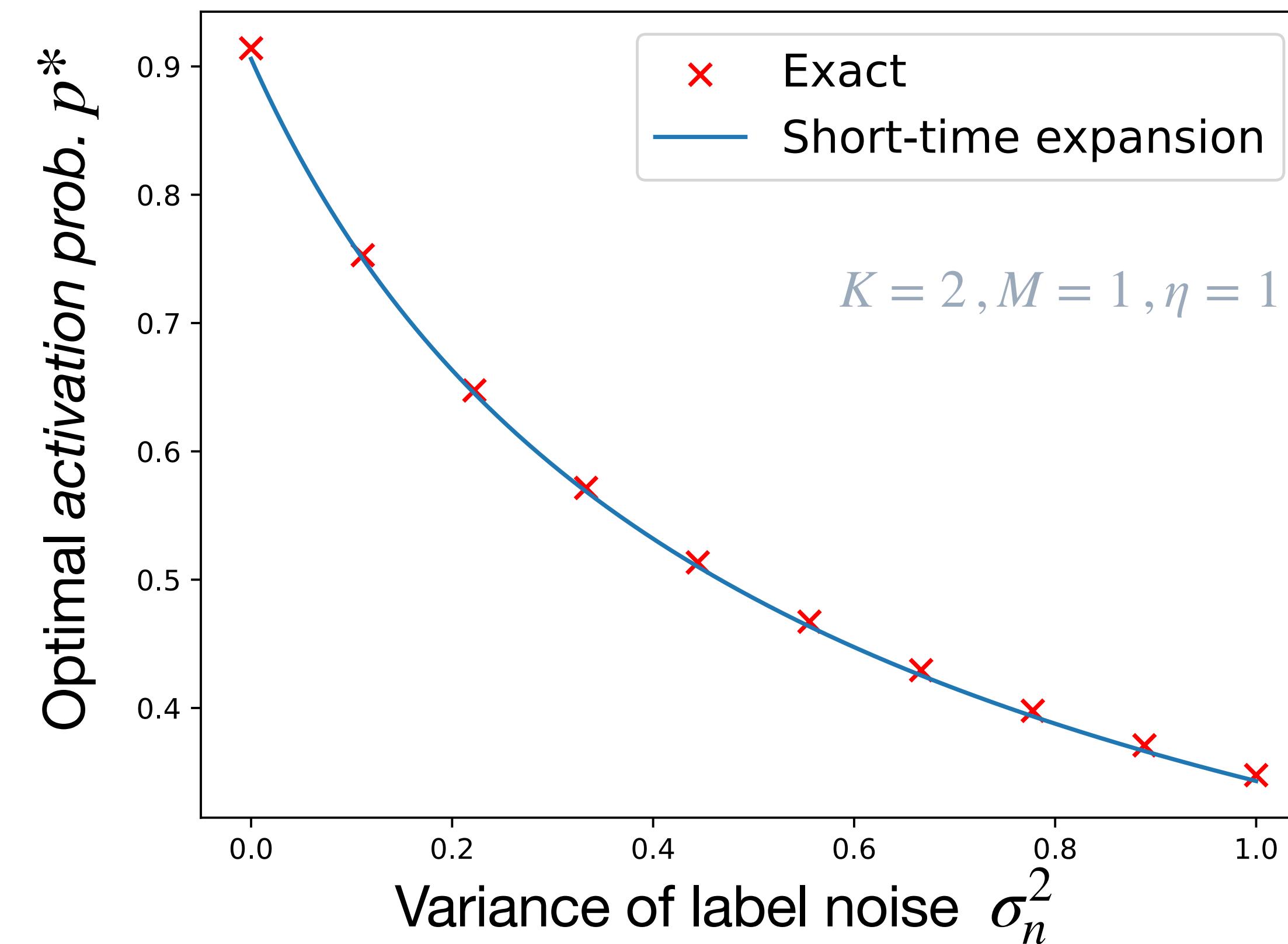
1. The optimal *constant* activation probability depends on label noise, width, and learning rate.
-

# Online learning with constant node-activation prob.

1. The optimal *constant* activation probability depends on label noise, width, and learning rate.

Early-time dynamics:

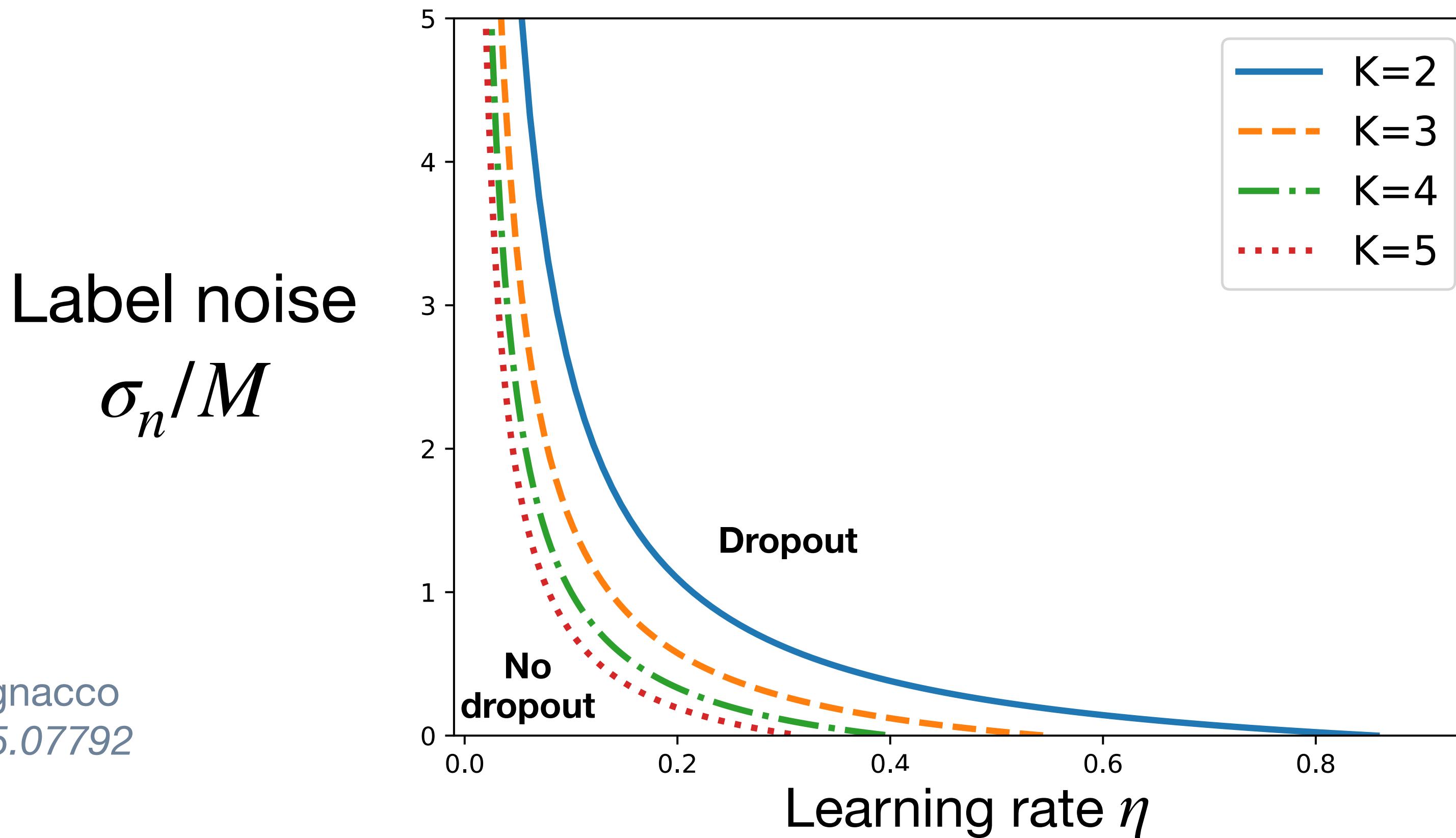
$$p^* = \min\left\{1; \frac{12M}{3M\eta + 9\eta\sigma_n^2 + \sqrt{3\eta(M + 3\sigma_n^2)[M(32(K - 1) + 3\eta) + 9\eta\sigma_n^2]}}\right\}$$



# Online learning with constant node-activation prob.

1. The optimal *constant* activation probability depends on label noise, width, and learning rate.

Early-time dynamics:  $p^* = \min\{1; \frac{12M}{3M\eta + 9\eta\sigma_n^2 + \sqrt{3\eta(M + 3\sigma_n^2)[M(32(K - 1) + 3\eta) + 9\eta\sigma_n^2]}}\}$



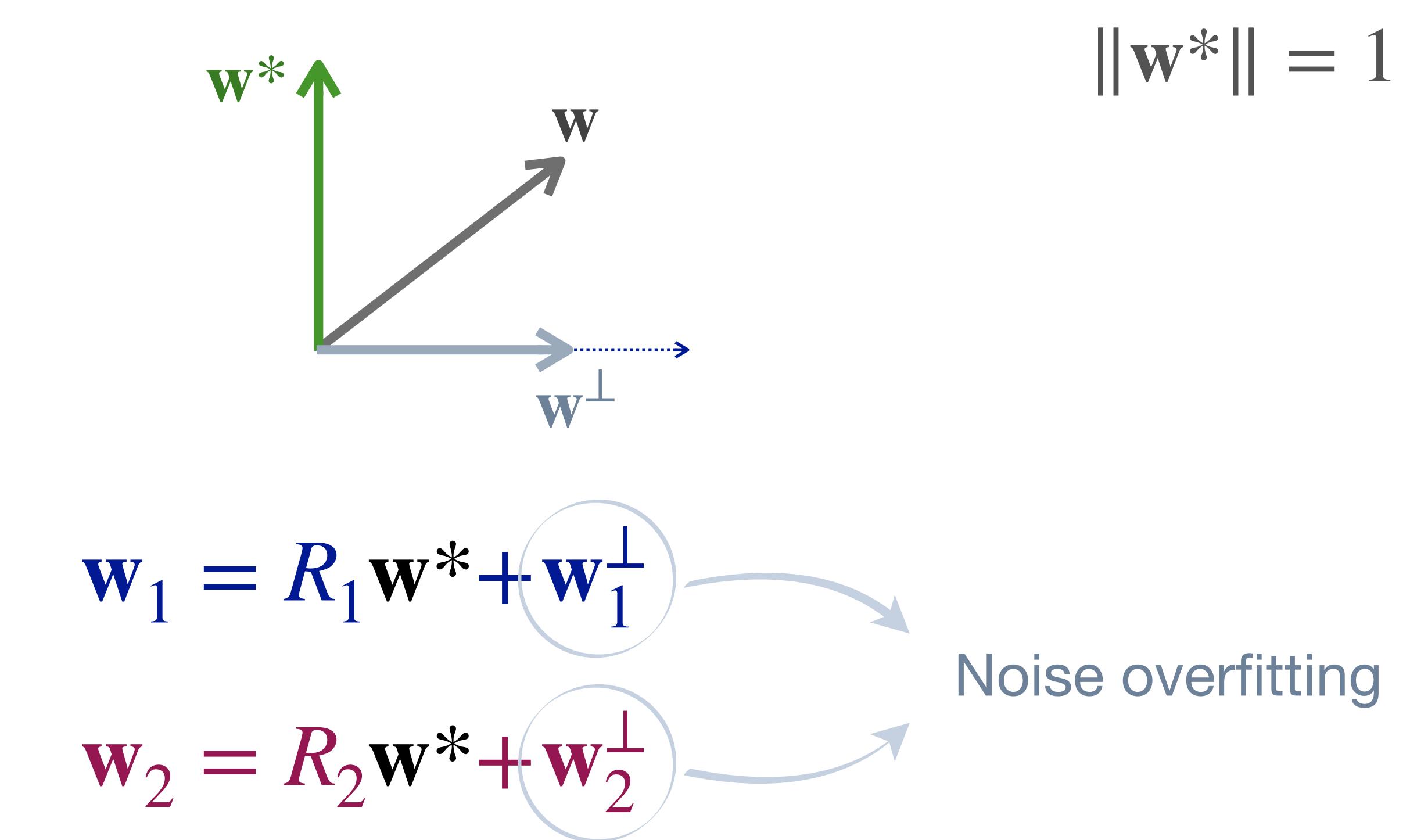
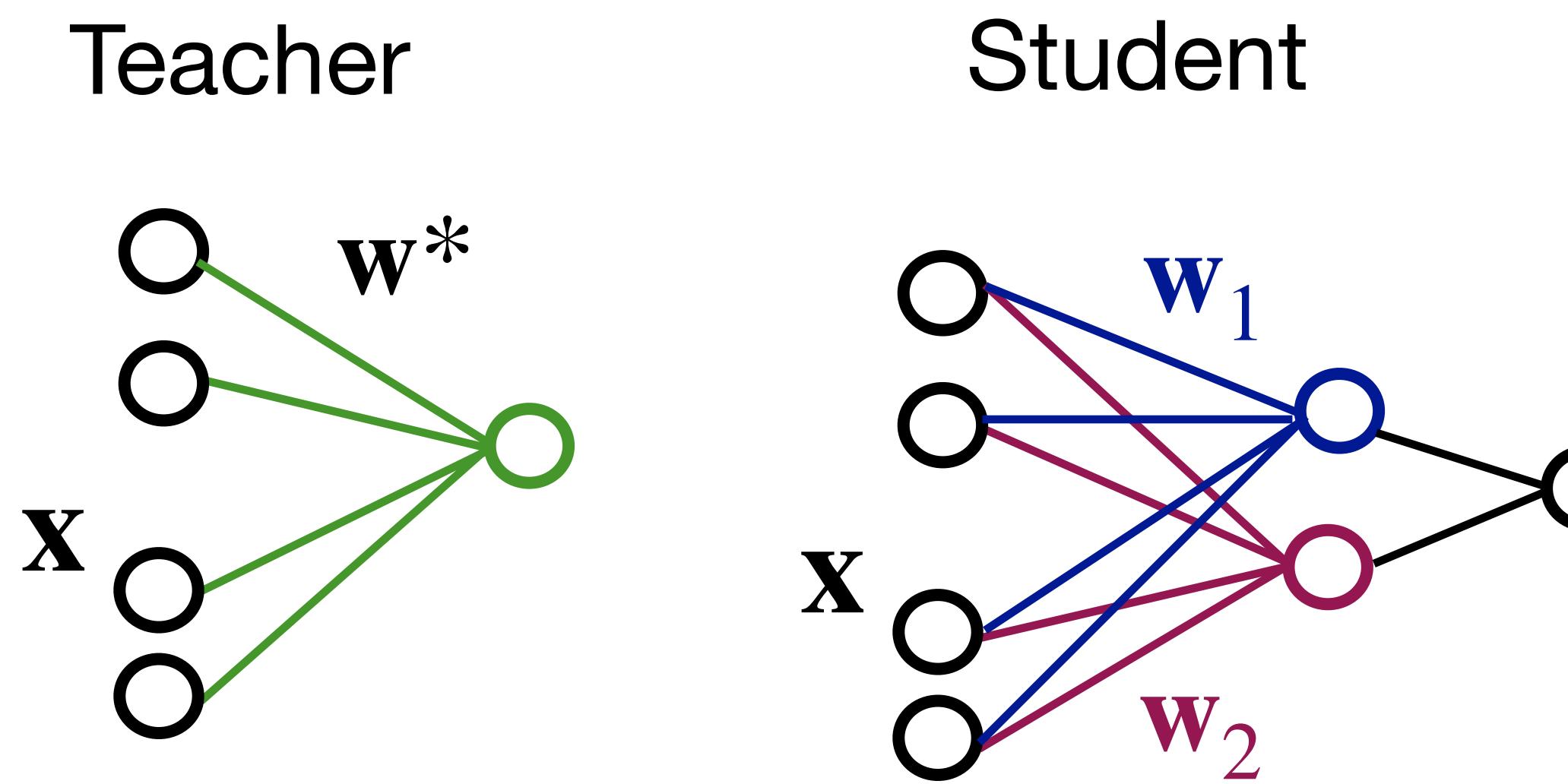
Dropout is optimal if:

$$\sigma_n^2 \geq M \left[ \frac{2}{(4K - 1)\eta} - \frac{1}{3} \right]$$

# Online learning with constant node-activation prob.

1. The optimal *constant* activation probability depends on label noise, width, and learning rate.
2. Dropout effectively decorrelates the hidden units of the student network.

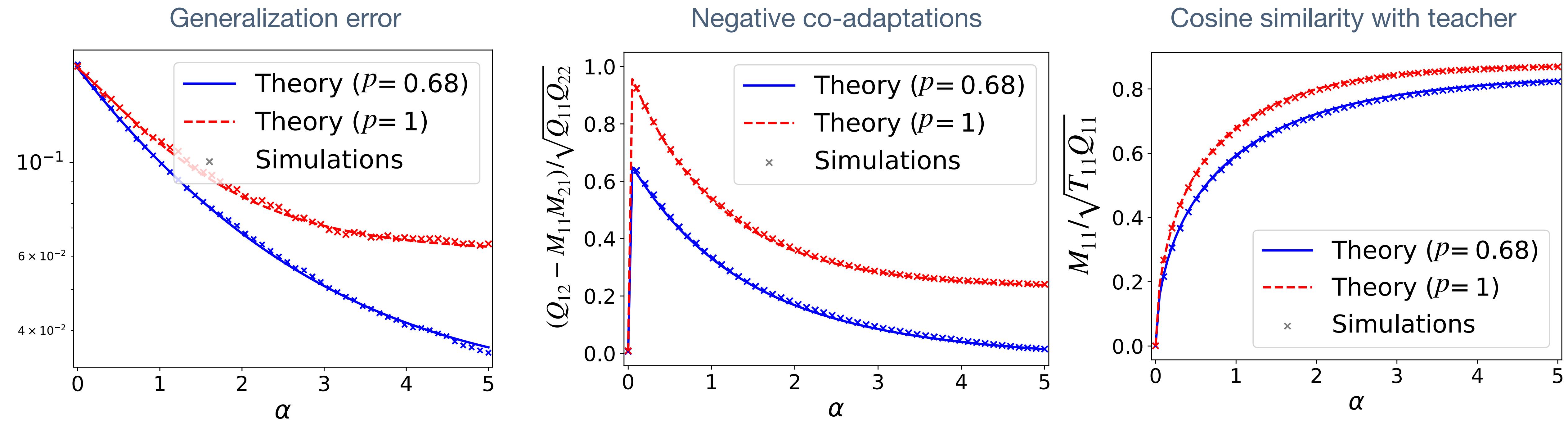
e.g., take  $M = 1$  and  $K = 2$



# Online learning with constant node-activation prob.

1. The optimal *constant* activation probability depends on label noise, width, and learning rate.
2. Dropout effectively decorrelates the hidden units of the student network.

e.g., take  $M = 1$  and  $K = 2$



# *Adaptive dropout probability*

# ANNEALED DROPOUT TRAINING OF DEEP NETWORKS

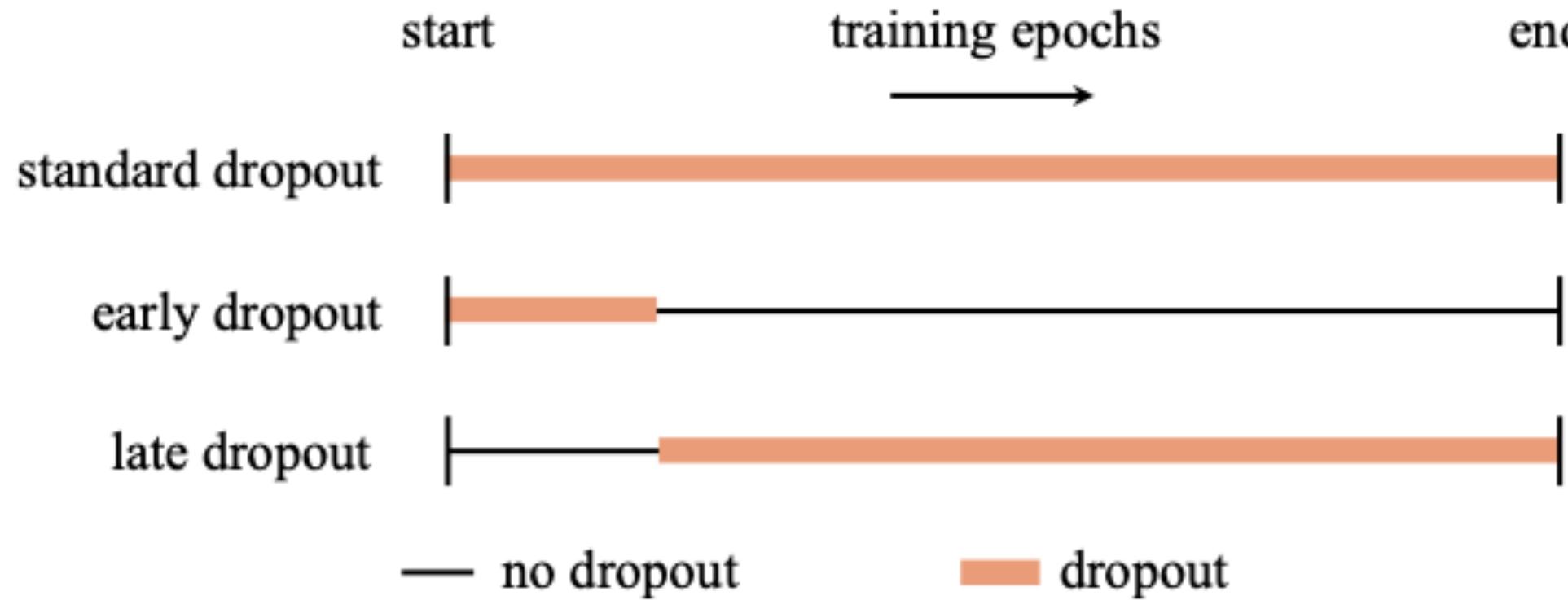
*Steven J. Rennie, Vaibhava Goel, and Samuel Thoma*

IBM Thomas J. Watson Research Center  
NY, USA

{sjrennie, vgoel, sthomas}@us.ibm.com

## Dropout Reduces Underfitting

**Zhuang Liu<sup>\* 1</sup> Zhiqiu Xu<sup>\* 2</sup> Joseph Jin<sup>2</sup> Zhiqiang Shen<sup>3</sup> Trevor Darrel**



## **Curriculum Dropout**

Pietro Morerio<sup>1</sup>, Jacopo Cavazza<sup>1,2</sup>, Riccardo Volpi<sup>1,2</sup>, René Vidal<sup>3</sup> and Vittorio Murino<sup>1,4</sup>

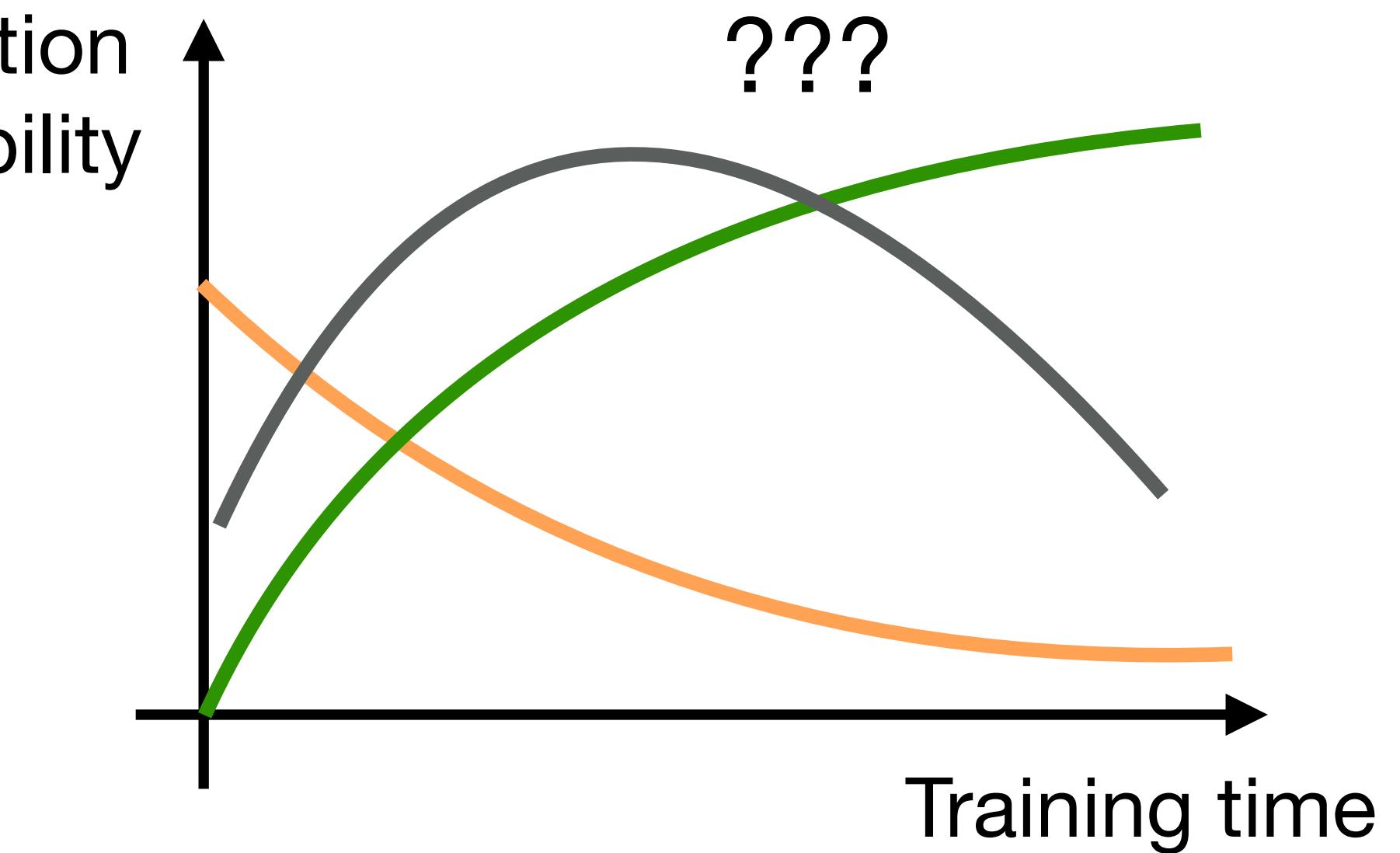
<sup>1</sup>Pattern Analysis & Computer Vision (PAVIS) – Istituto Italiano di Tecnologia – Genova, 16163, Italy

<sup>2</sup>Electrical, Electronics and Telecommunication Engineering and Naval Architecture Department (DITEN) – Università degli Studi di Genova – *Genova, 16145, Italy*

<sup>3</sup>Department of Biomedical Engineering – Johns Hopkins University – Baltimore, MD 21218, USA

<sup>4</sup>Computer Science Department – Università di Verona – Verona, 37134, Italy

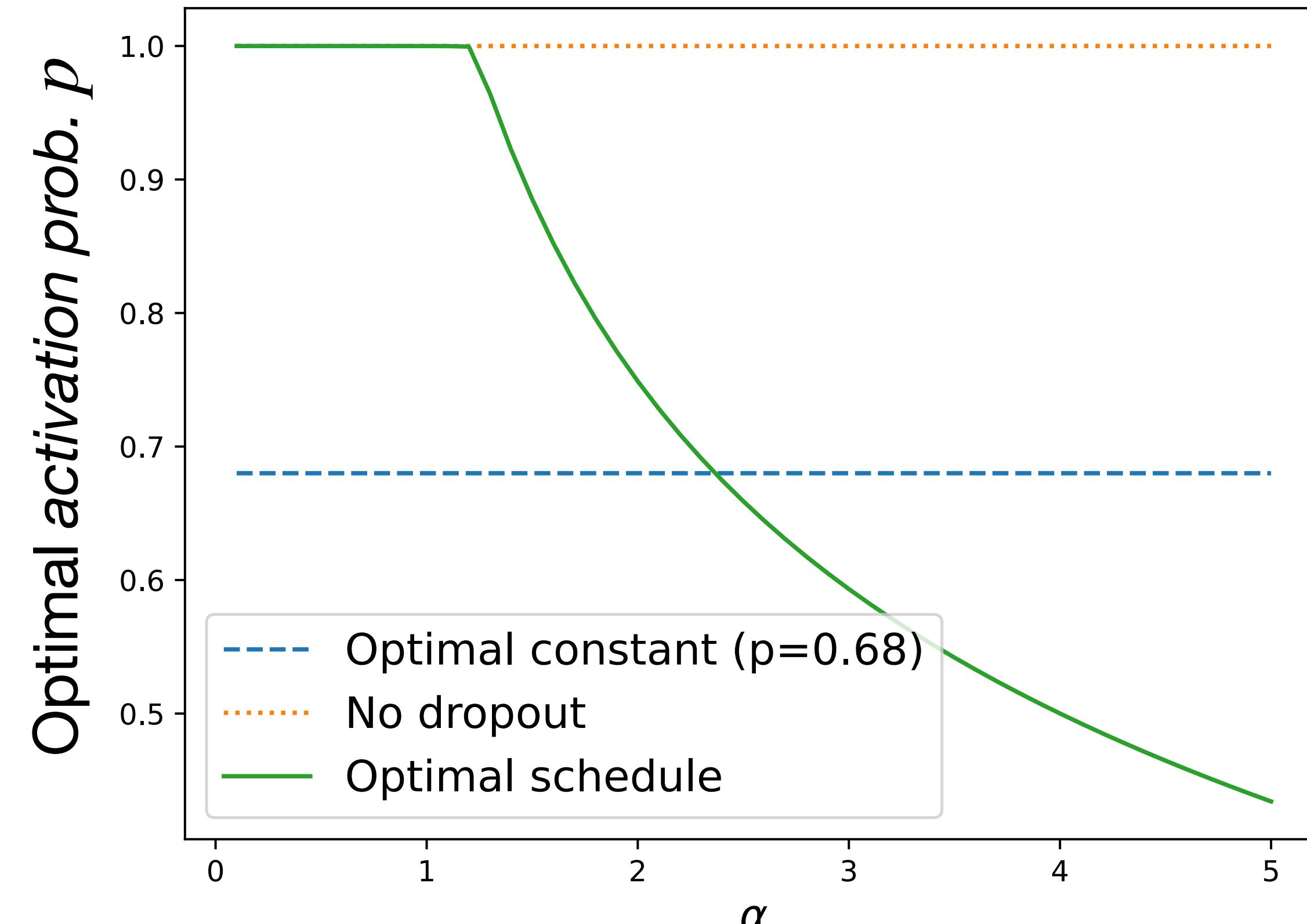
{pietro.morerio, jacopo.cavazza, riccardo.volpi, vittorio.murino}@iit.it, rvidal@cis.jhu.edu



# Optimal dropout schedules

Teacher-student setting ( $K = 2, M = 1$ )

1. The optimal schedule monotonically decreases the node-activation probability.



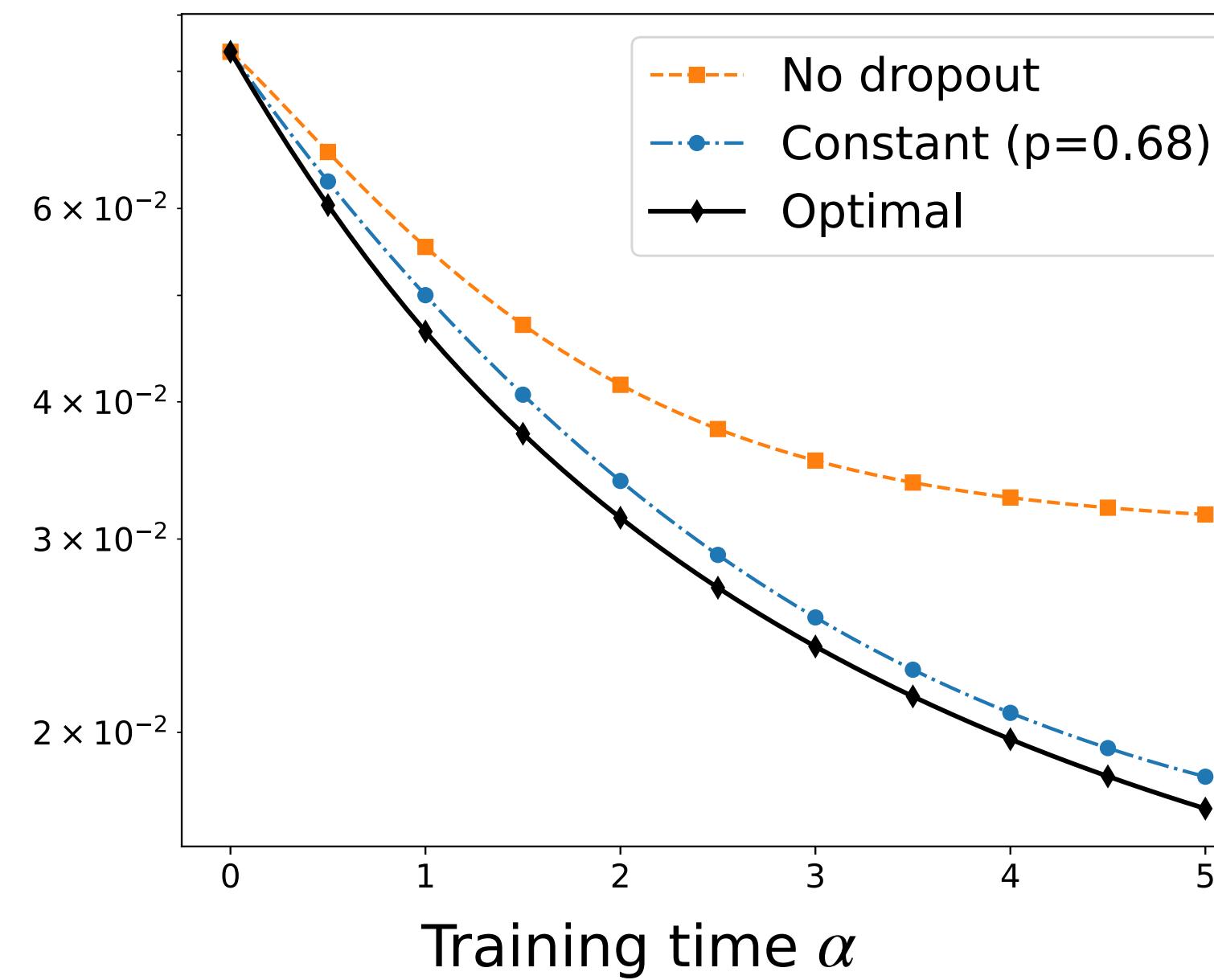
$$\eta = 1, \alpha_F = 5, \sigma_n = 0.3, Q_0 = 0$$

# Optimal dropout schedules

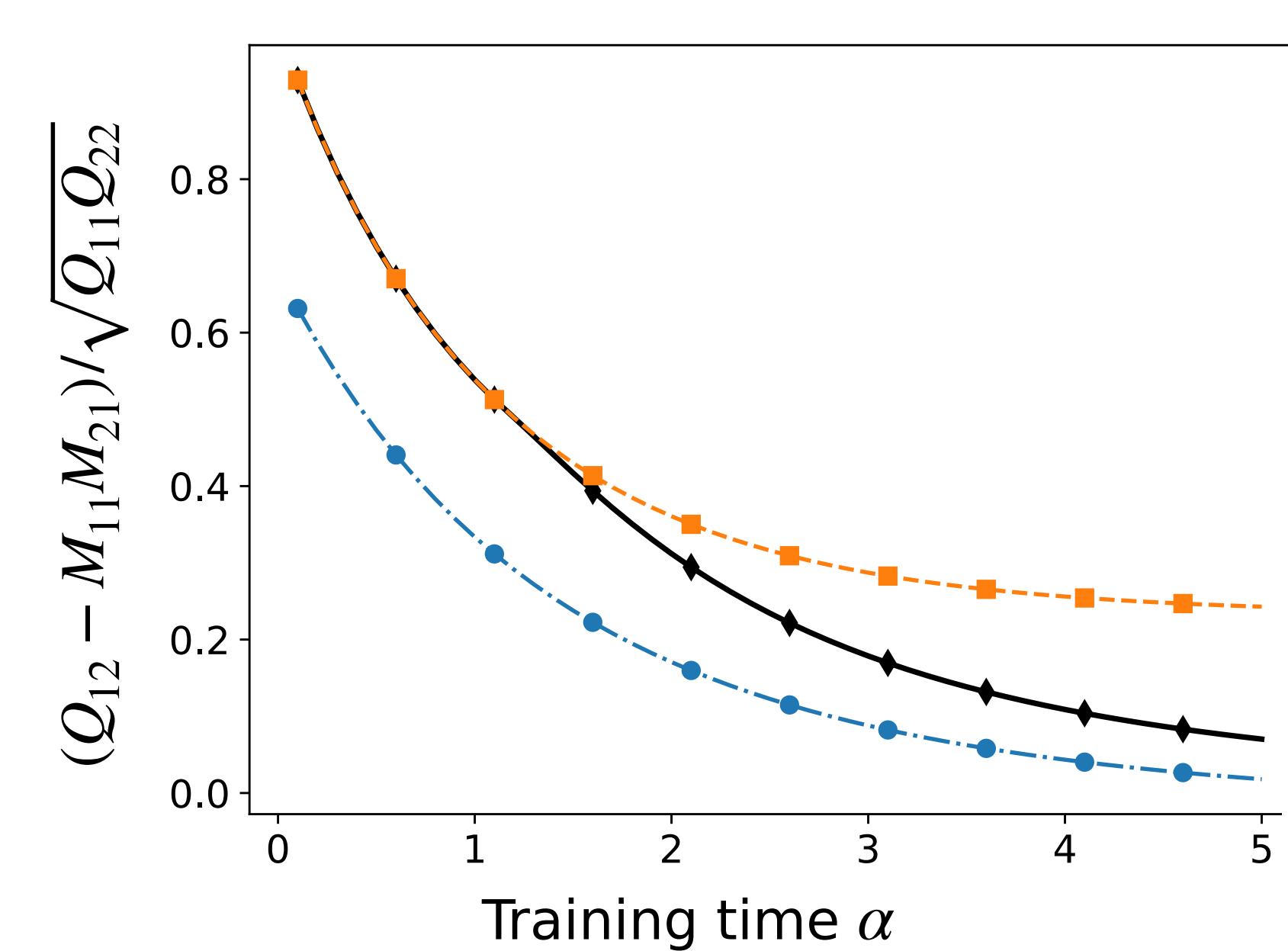
Teacher-student setting ( $K = 2, M = 1$ )

1. The optimal schedule monotonically decreases the node-activation probability.
2. Adaptive dropout achieves the optimal balance between node decorrelation and signal recovery.

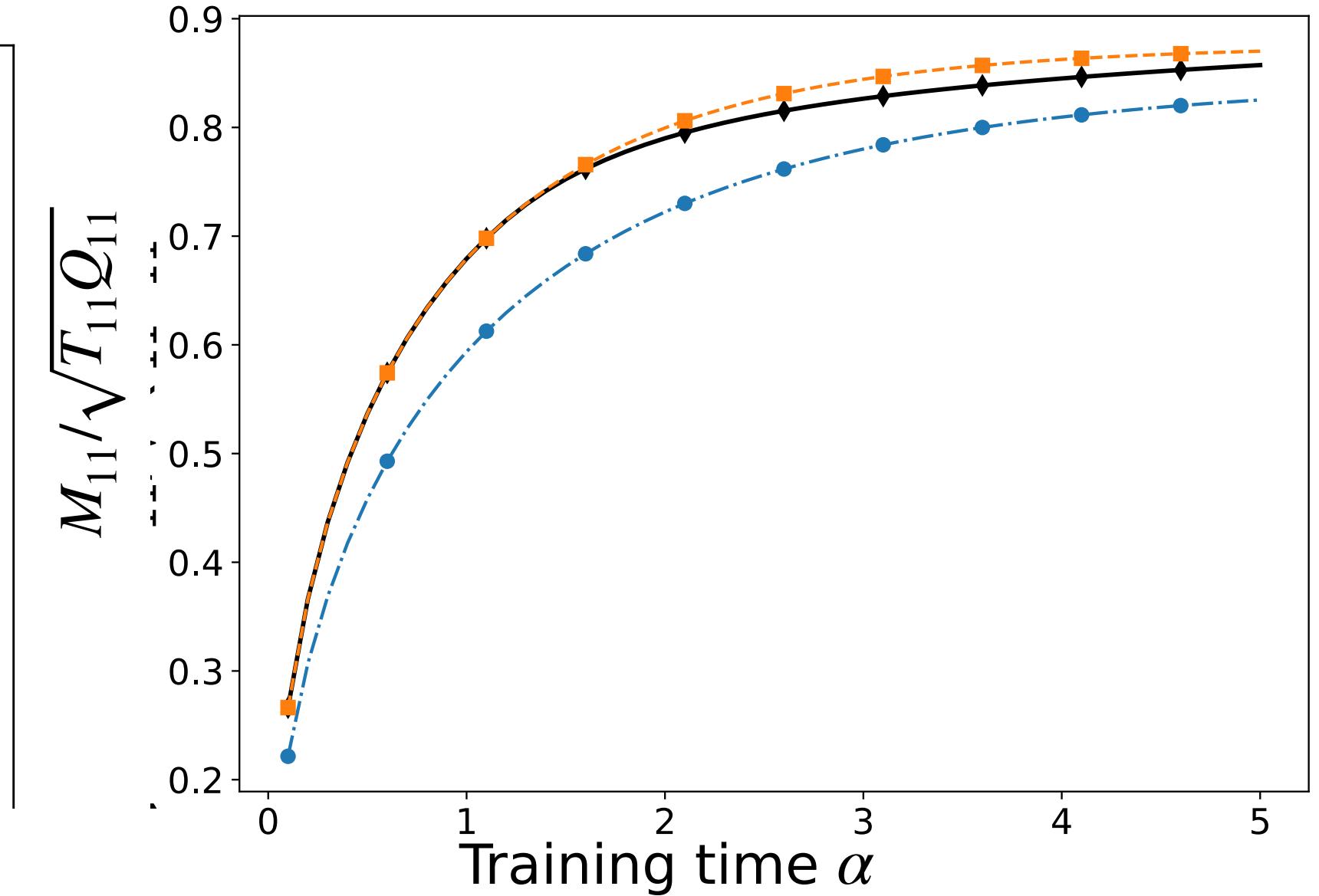
Generalization error



Negative co-adaptations



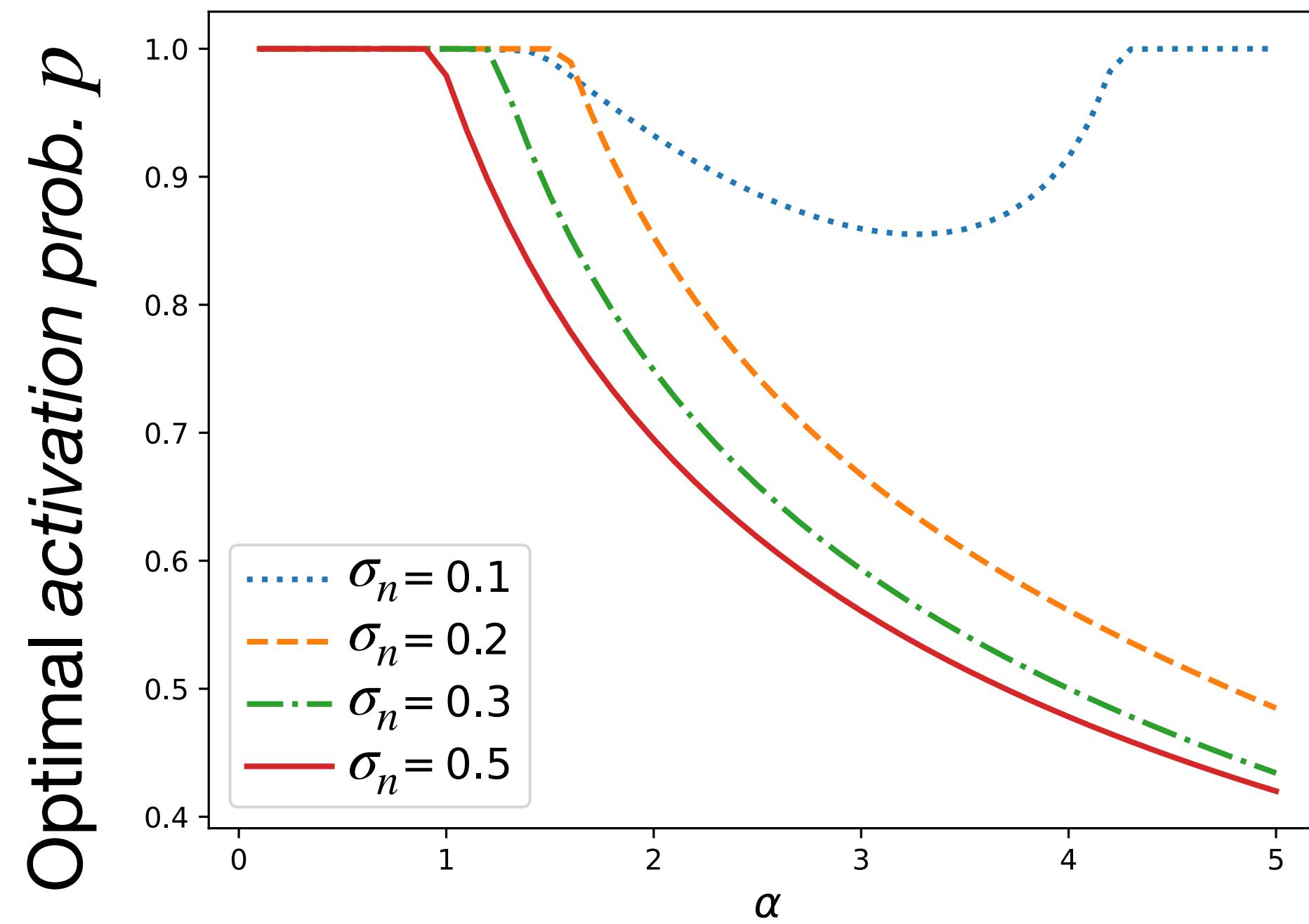
Cosine similarity with teacher



# Optimal dropout schedules

Teacher-student setting ( $K = 2, M = 1$ )

1. The optimal schedule monotonically decreases the node-activation probability.
2. Adaptive dropout achieves the optimal balance between node decorrelation and signal recovery.
3. The dropout schedule adapts to the noise level in the task.

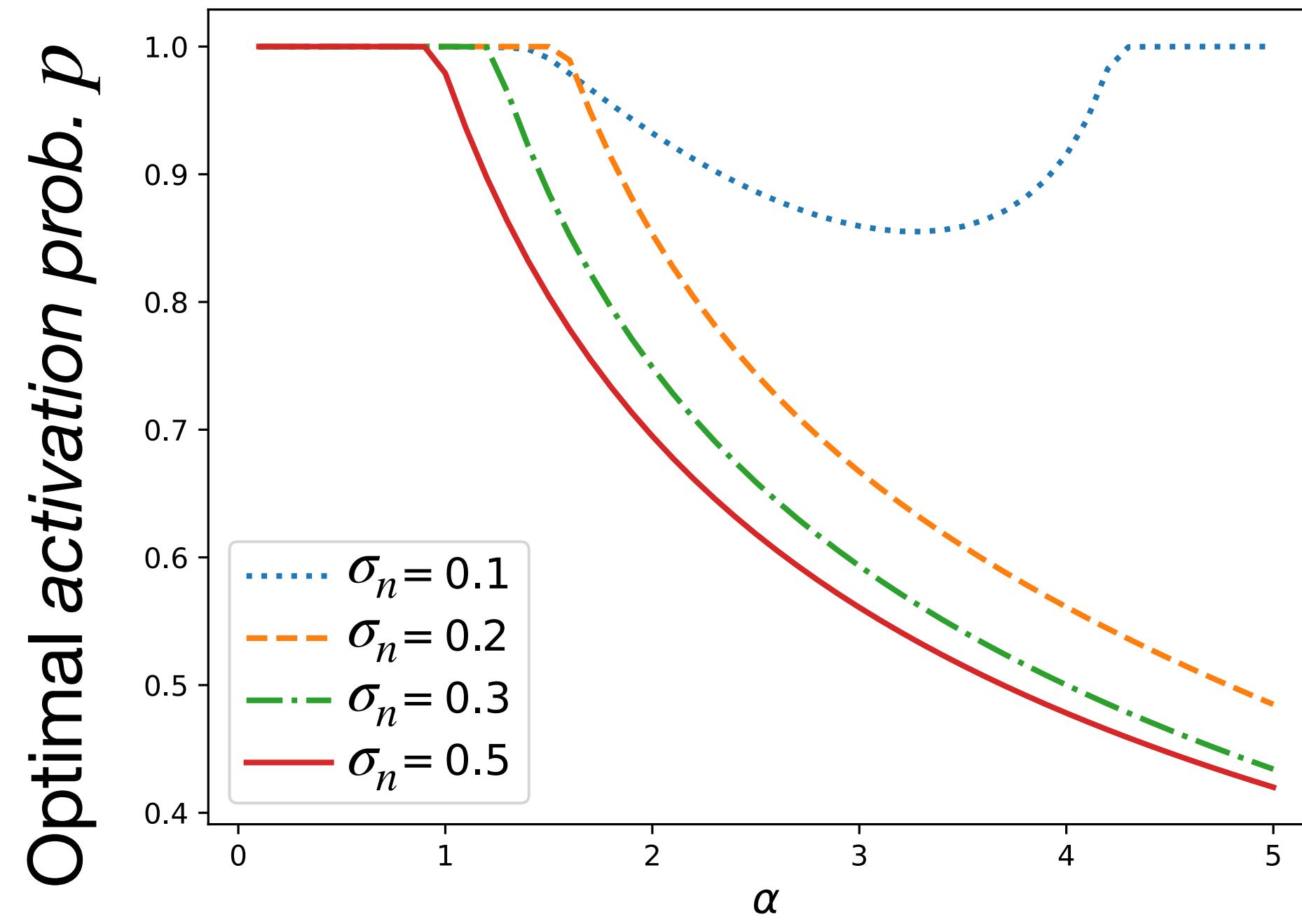


$$\eta = 1, \alpha_F = 5, Q_0 = 0$$

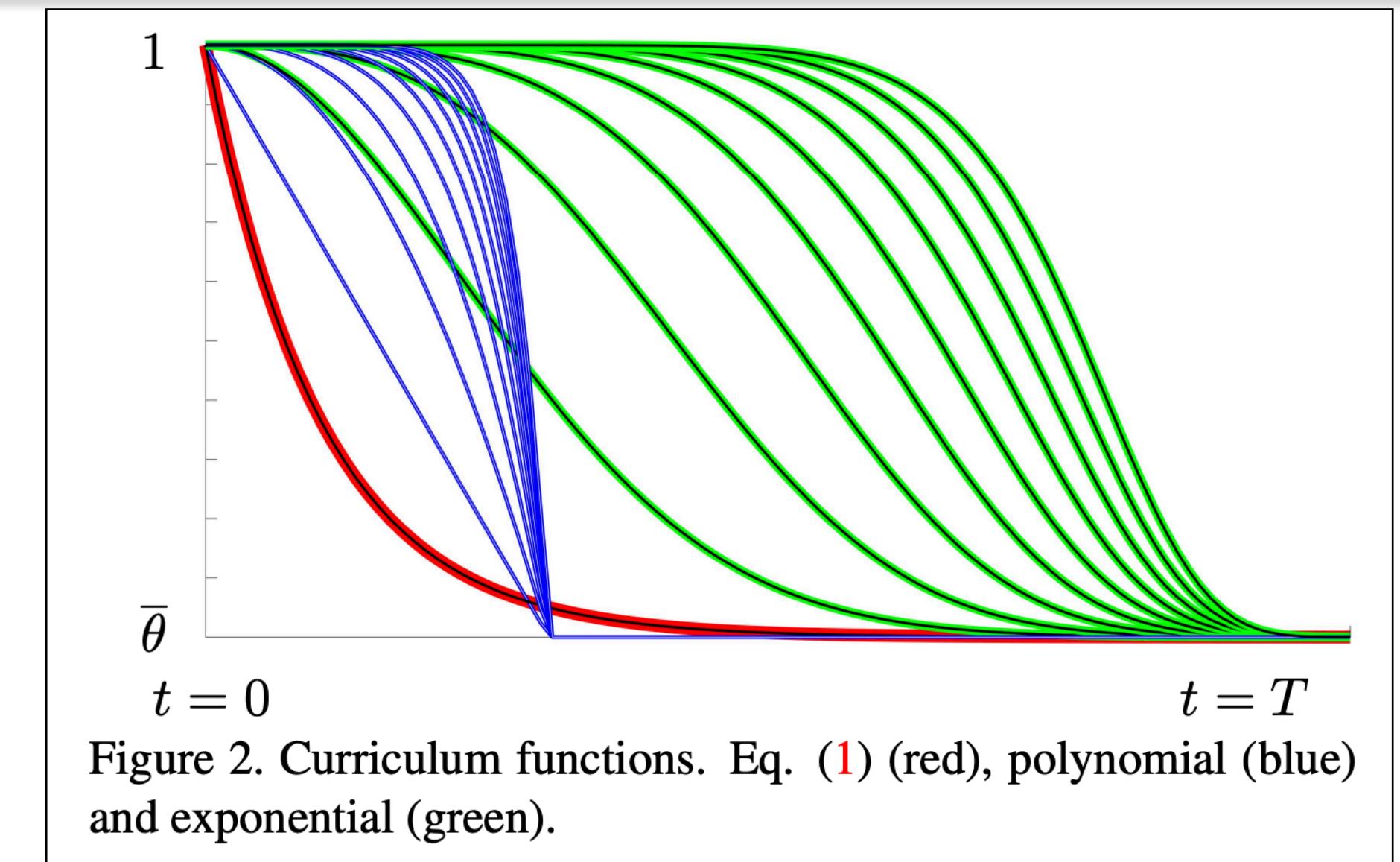
# Optimal dropout schedules

Teacher-student setting ( $K = 2, M = 1$ )

1. The optimal schedule monotonically decreases the node-activation rate.
2. Adaptive dropout achieves the optimal balance between node decorrelation and signal recovery.
3. The dropout schedule adapts to the noise level in the task.



$$\eta = 1, \alpha_F = 5, Q_0 = 0$$



Morerio, P., Cavazza, J., Volpi, R., Vidal, R., & Murino, V. (2017).  
*Curriculum dropout.*  
 In Proceedings of the IEEE international conference on computer vision (pp.  
 3544-3552)

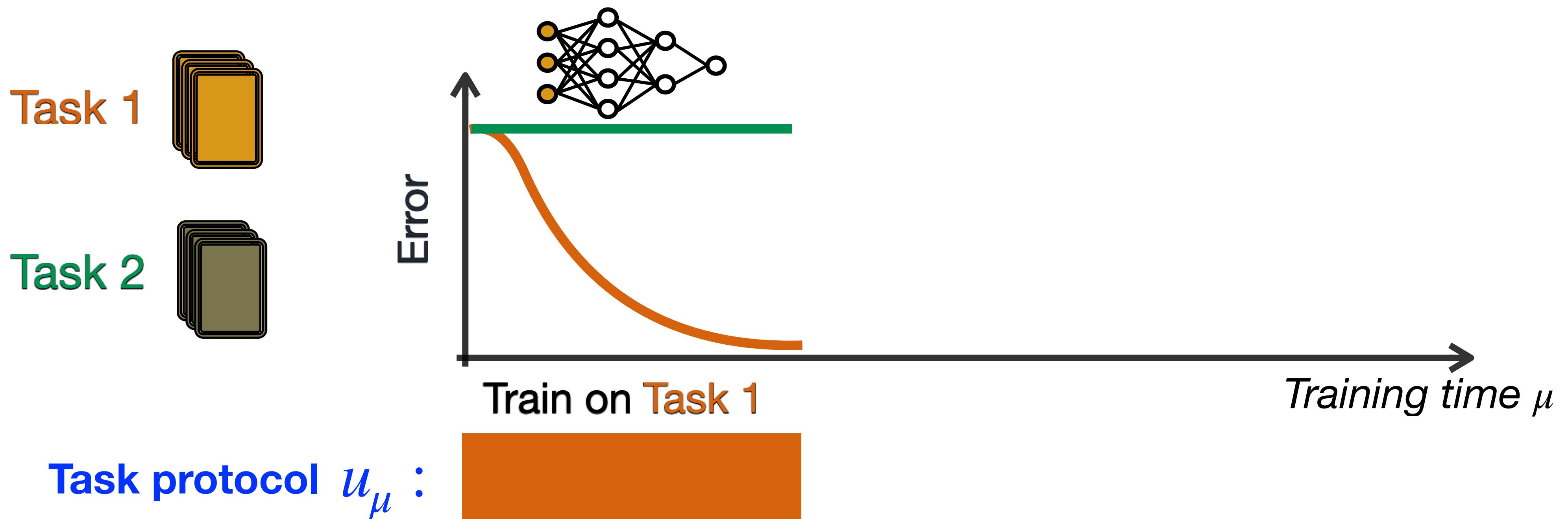
# Part II: *Continual learning*

# Continual learning

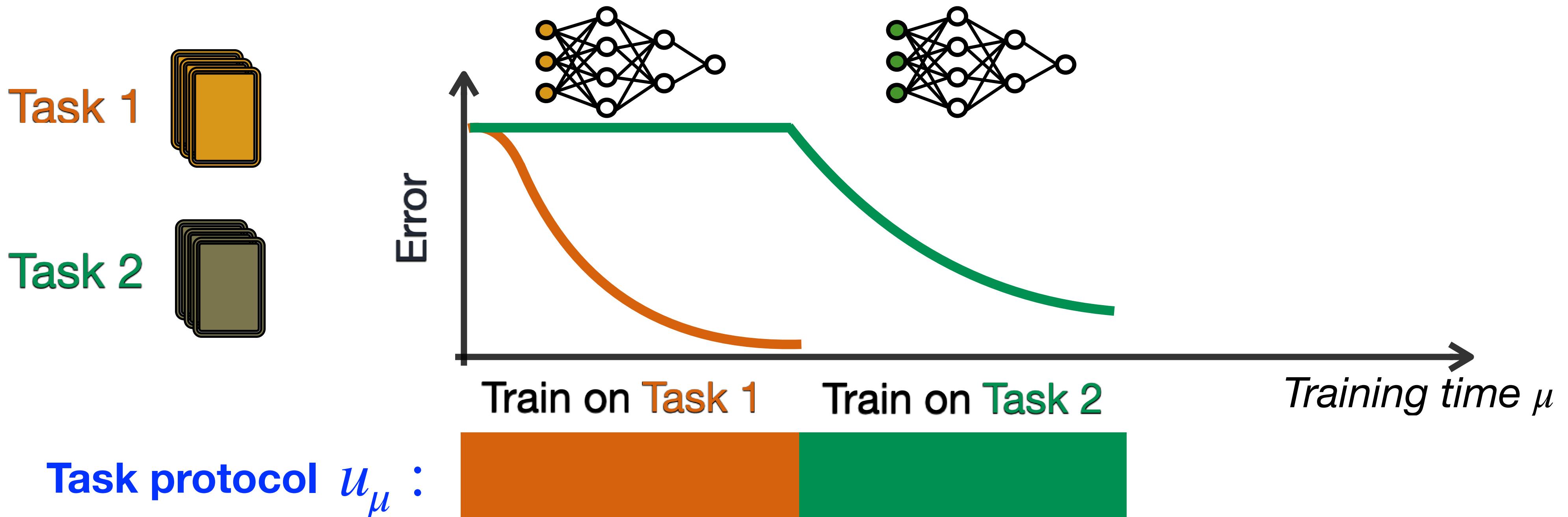


**Task protocol  $\mu_\mu$ :**

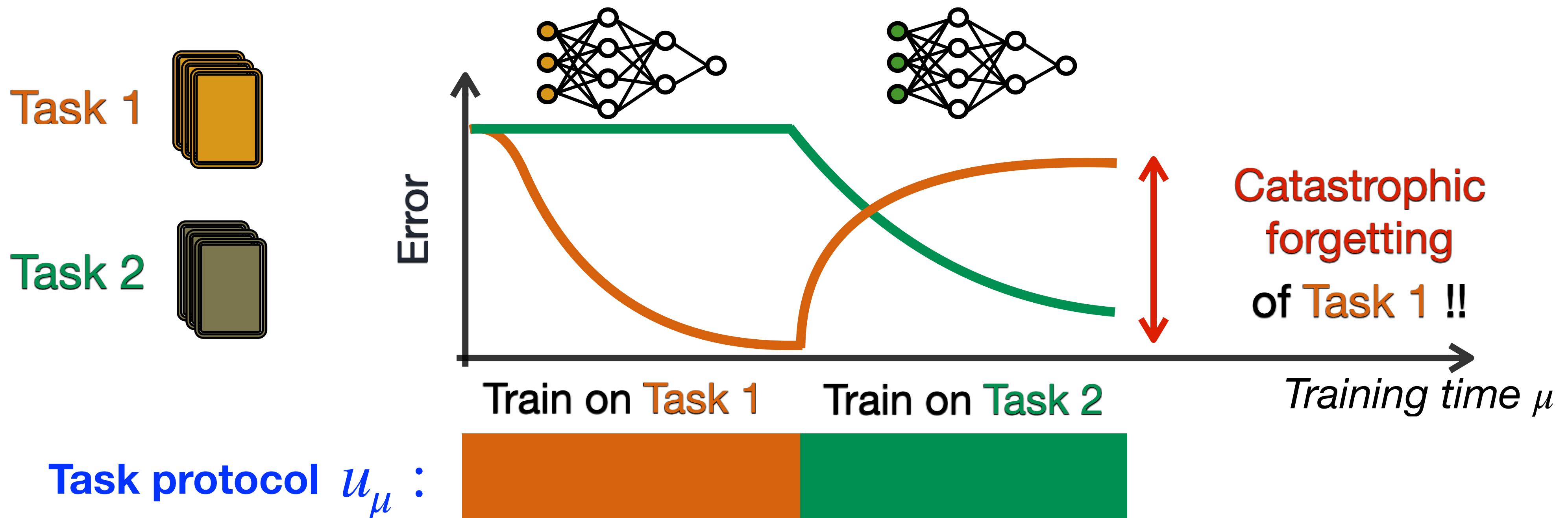
# Continual learning



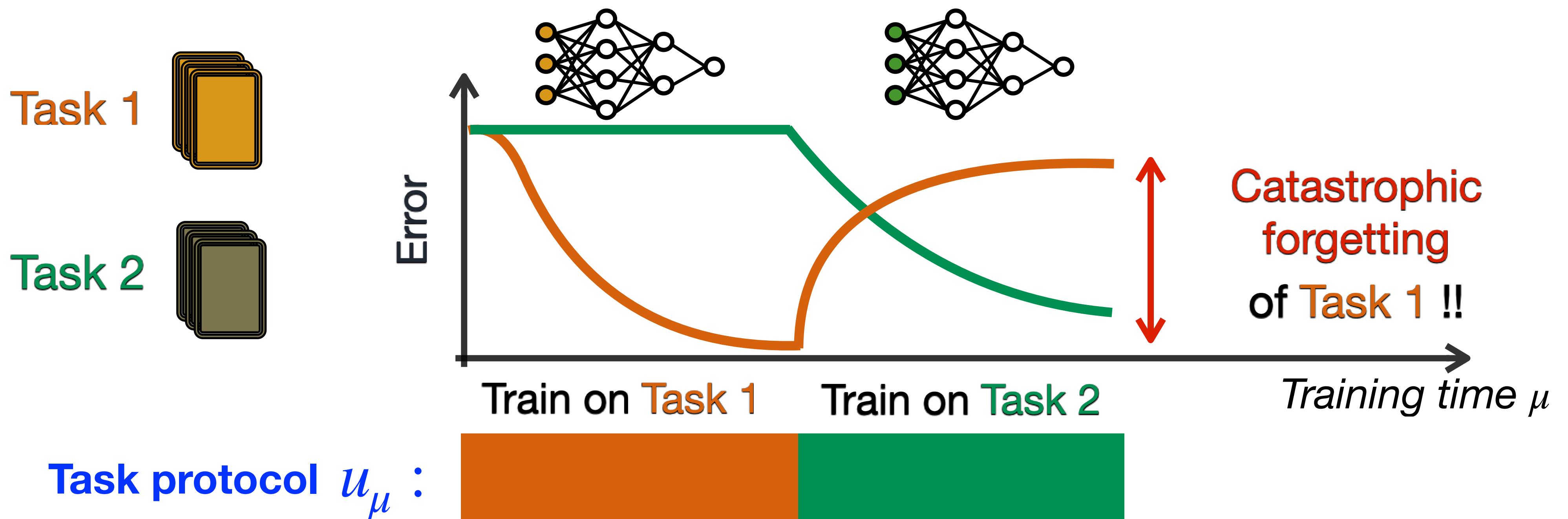
# Continual learning



# Continual learning



# Continual learning

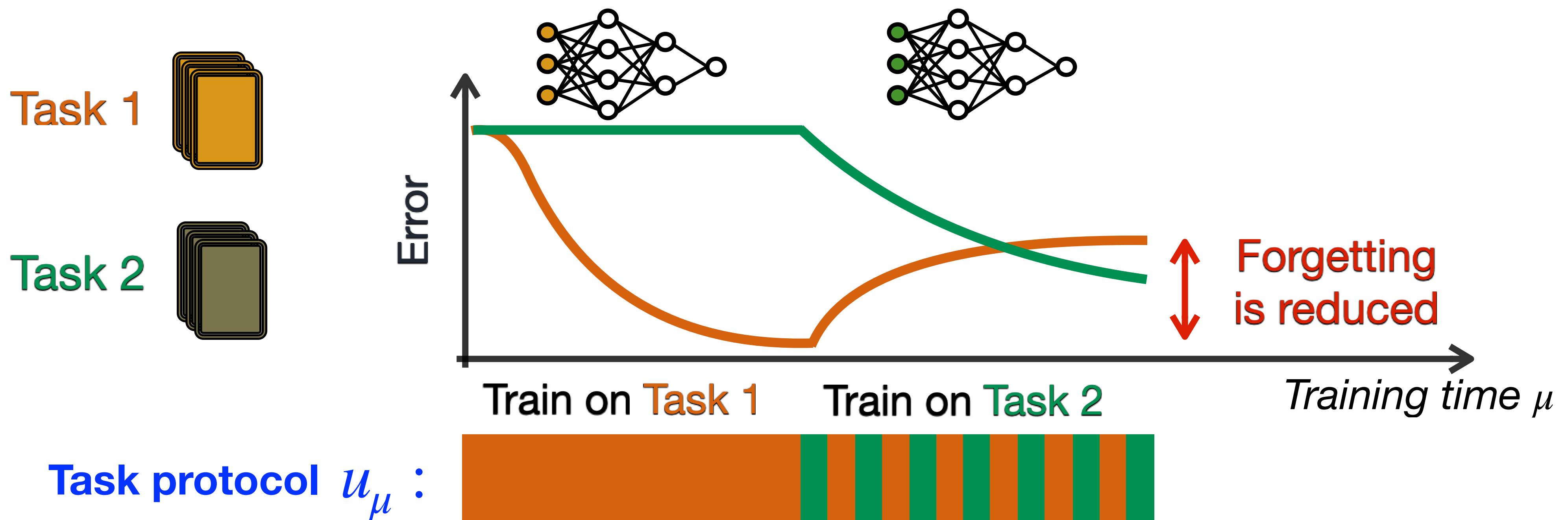


Humans & Animals

ML (empirical)

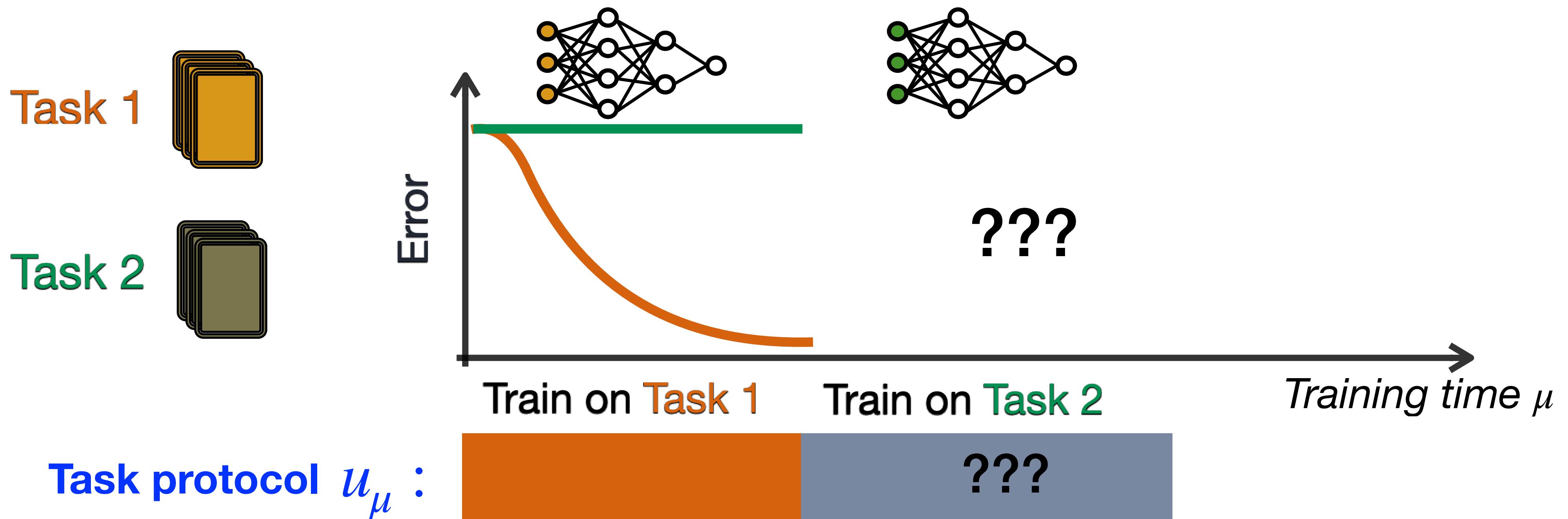
ML (theory)

# Continual learning



(Heuristic) interleaved *replay* strategy

# Continual learning



What is the *optimal* replay protocol?

# A teacher-student model of continual learning

Introduced in: Lee, Goldt, & Saxe (*ICML 2021*)

## Task 1

$$\mathcal{D}_1 = \{\mathbf{x}_\mu^{(1)}, y_\mu^{(1)}\}$$

$$\mathbf{x}_\mu^{(1)} \sim \mathcal{N}(0, I_N) \in \mathbf{R}^N$$

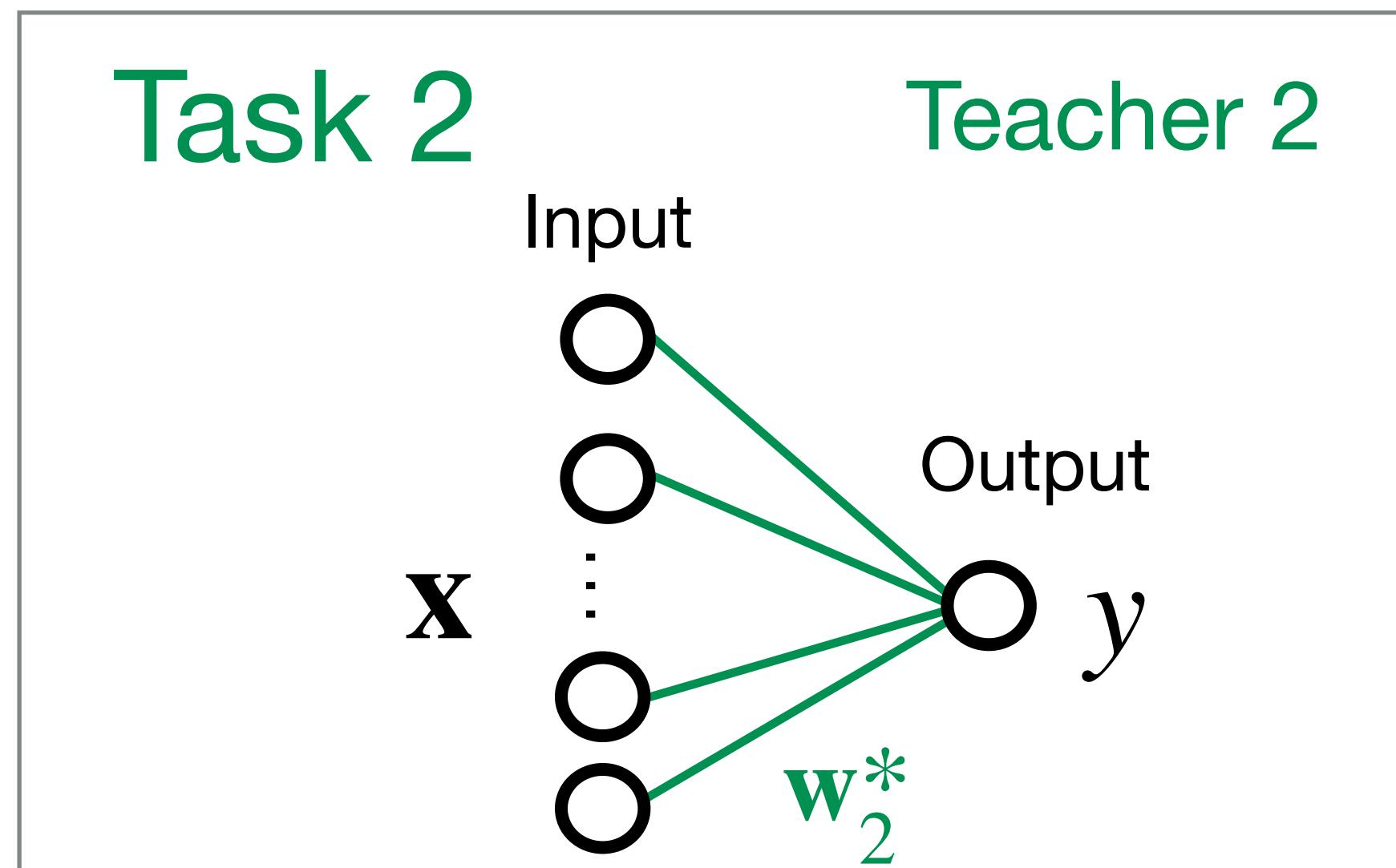
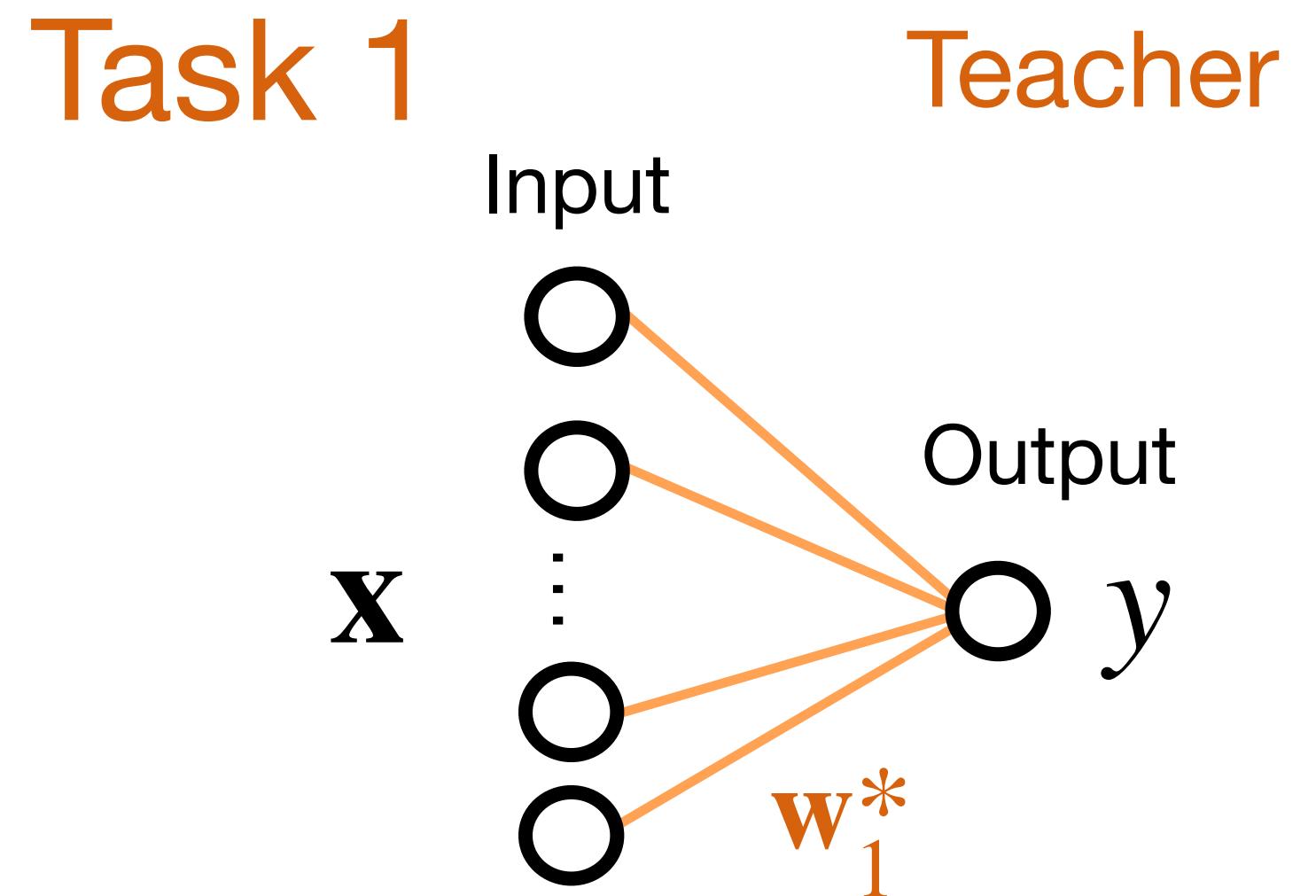
## Task 2

$$\mathcal{D}_2 = \{\mathbf{x}_\mu^{(2)}, y_\mu^{(2)}\}$$

$$\mathbf{x}_\mu^{(2)} \sim \mathcal{N}(0, I_N) \in \mathbf{R}^N$$

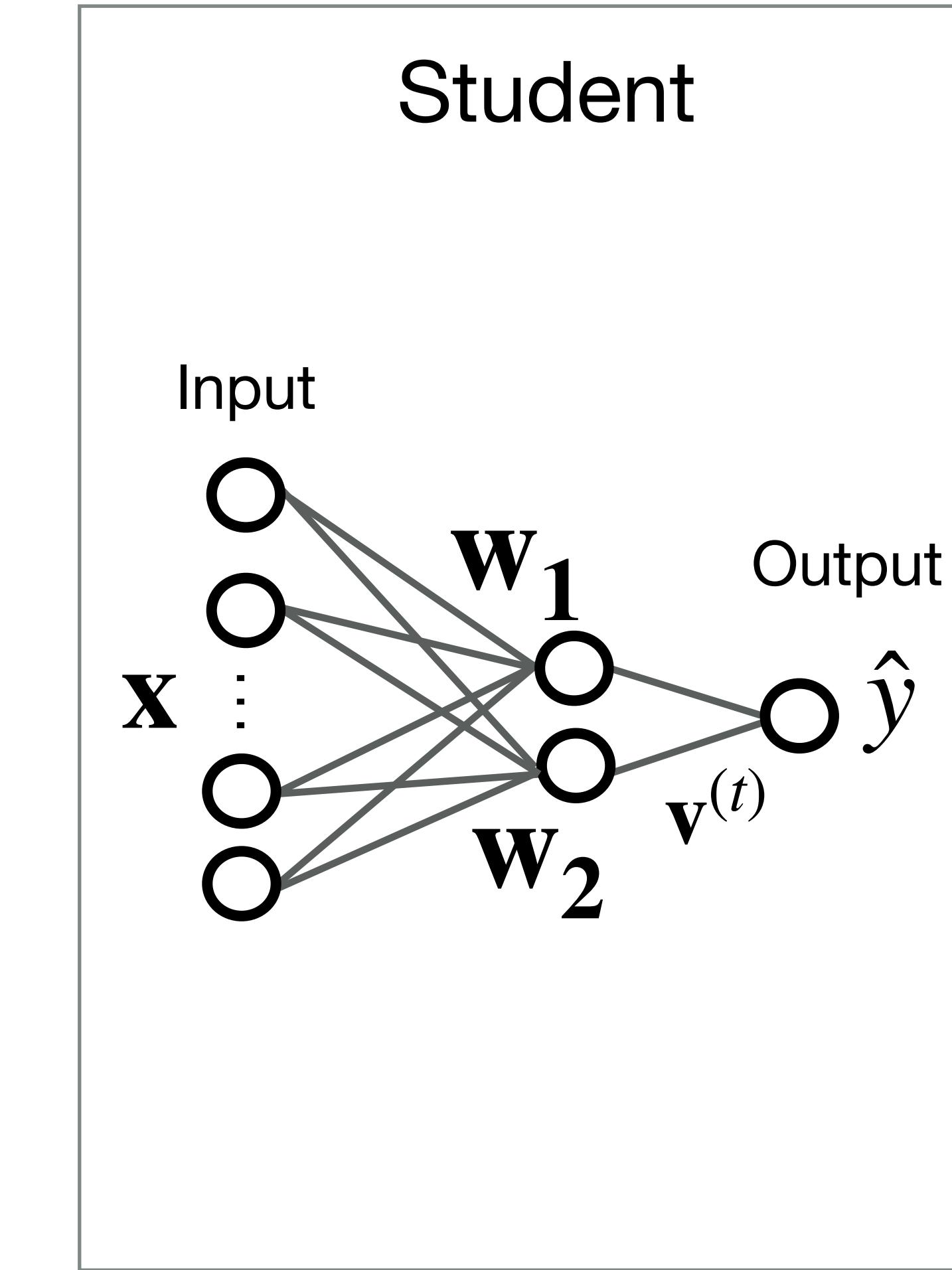
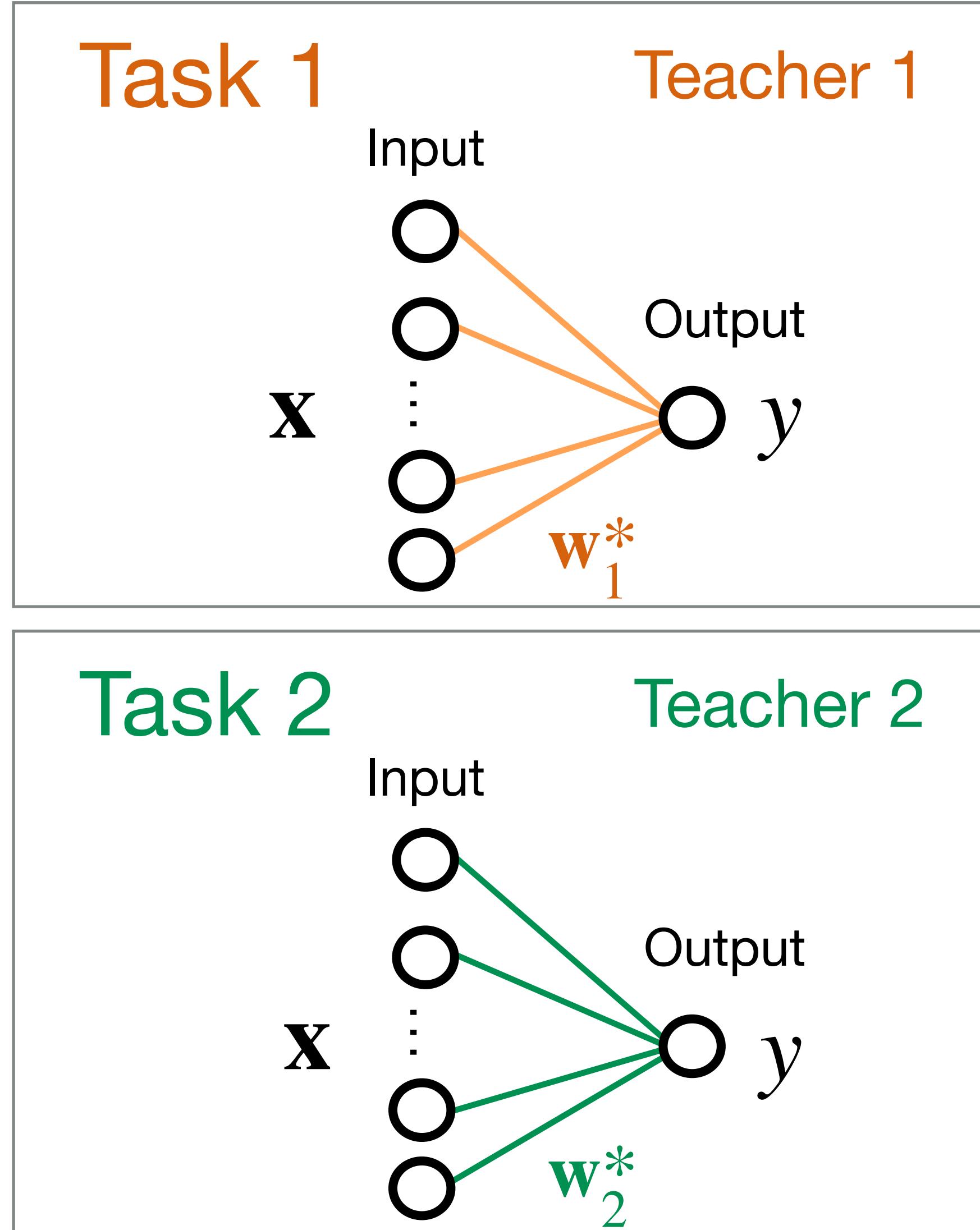
# A teacher-student model of continual learning

Introduced in: Lee, Goldt, & Saxe (ICML 2021)



# A teacher-student model of continual learning

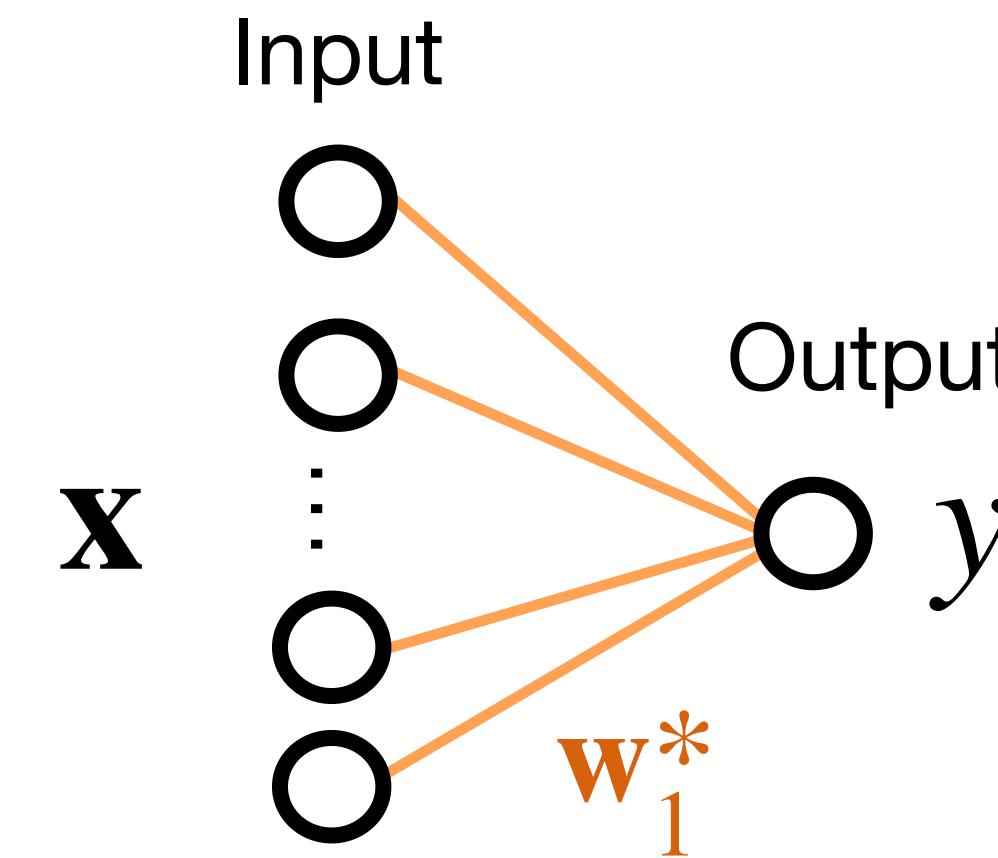
Introduced in: Lee, Goldt, & Saxe (ICML 2021)



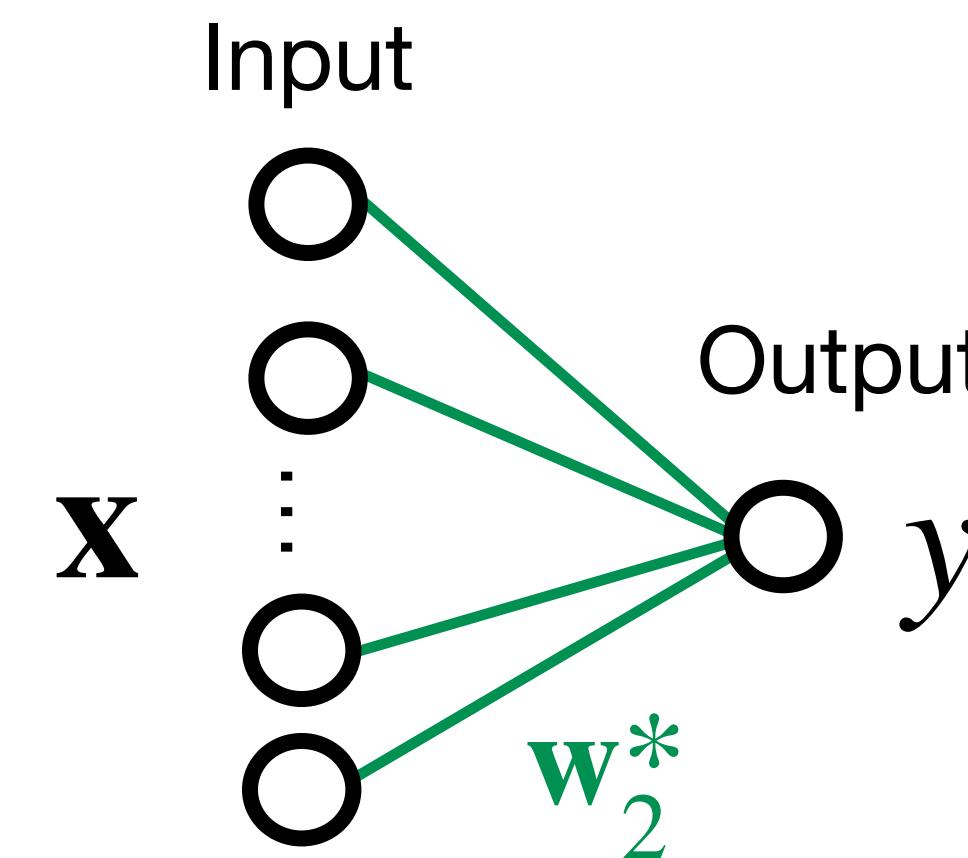
# A teacher-student model of continual learning

Introduced in: Lee, Goldt, & Saxe (ICML 2021)

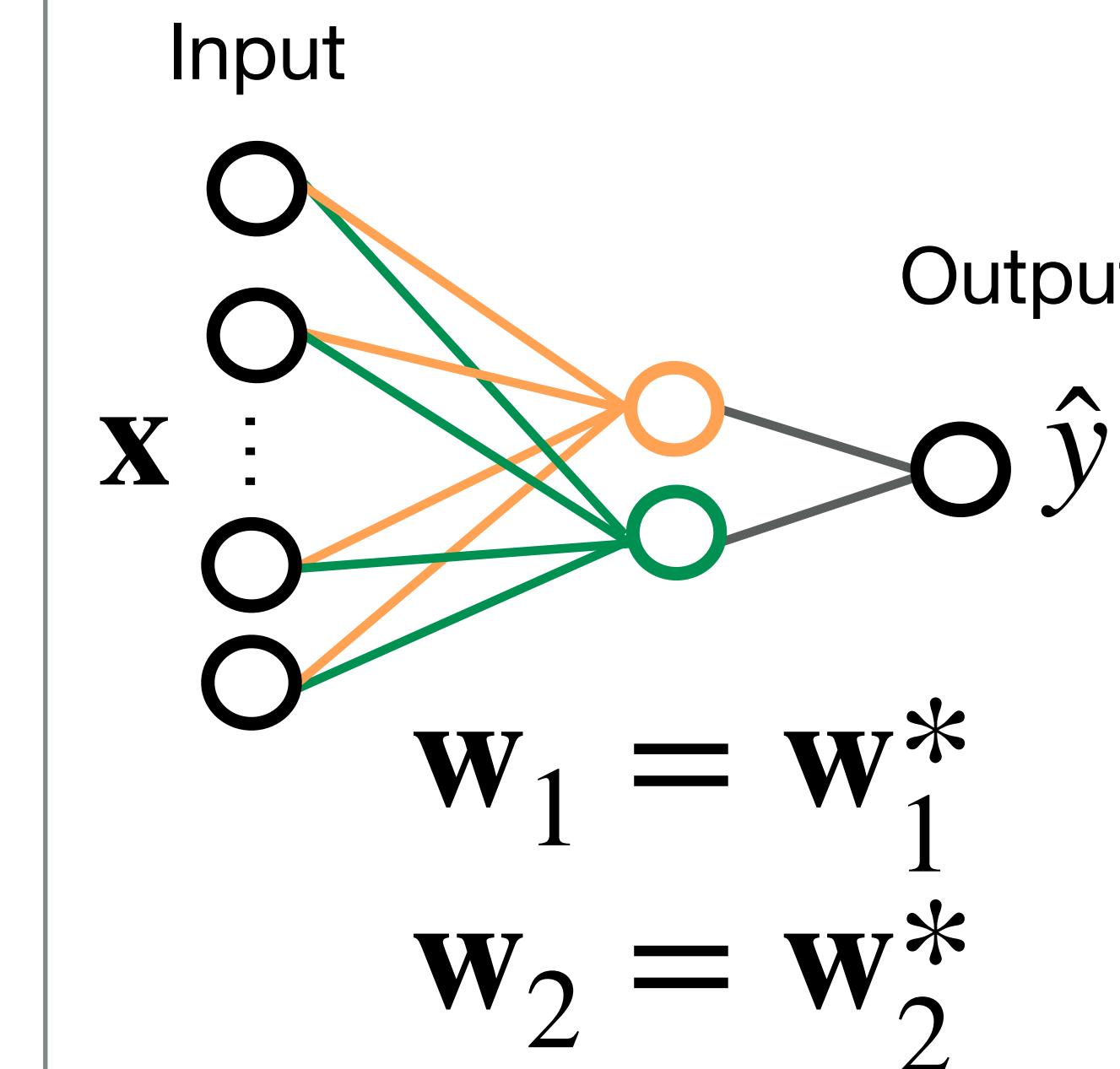
Task 1



Task 2



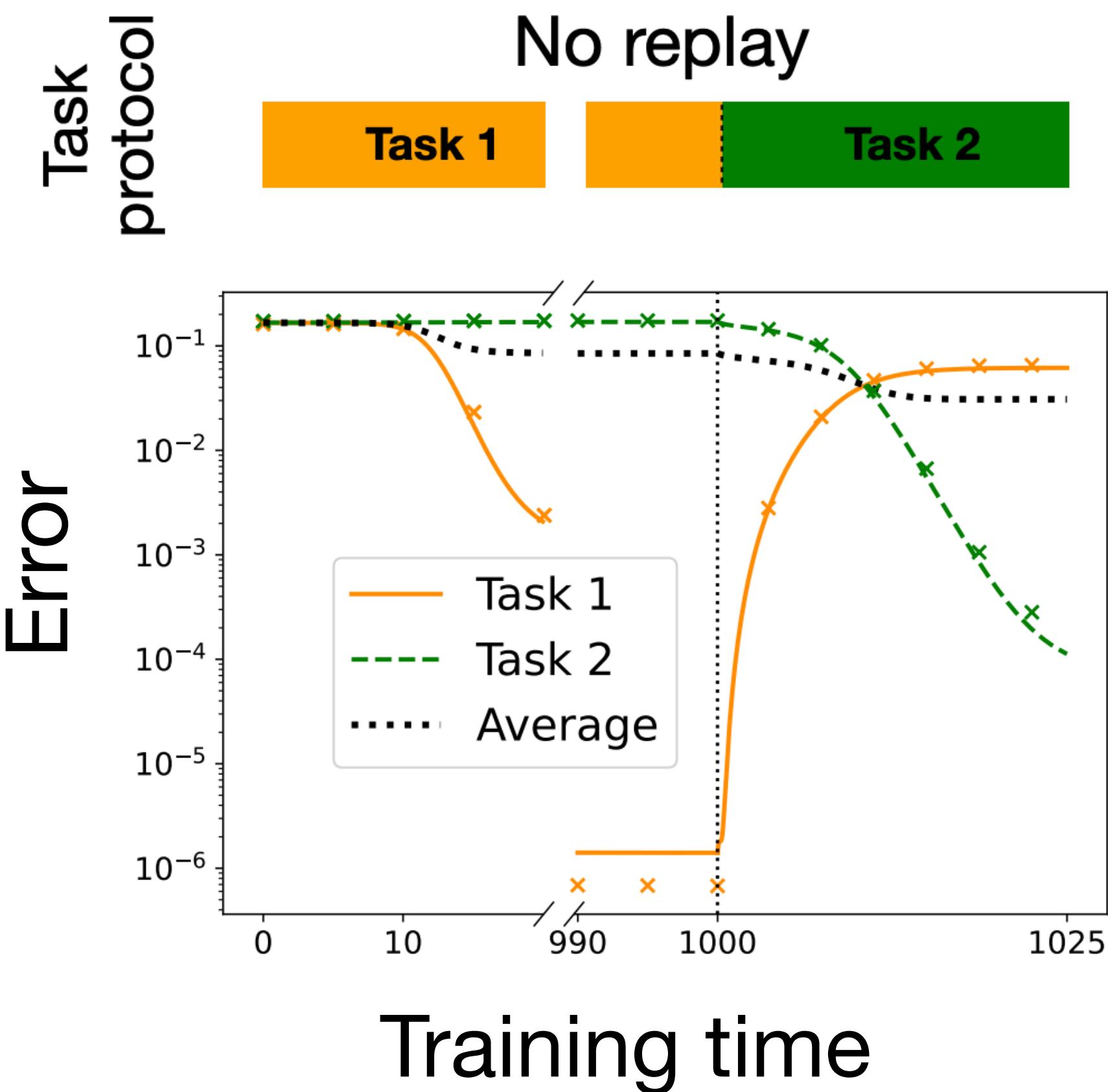
Student



The student has  
the capacity to  
learn both tasks.

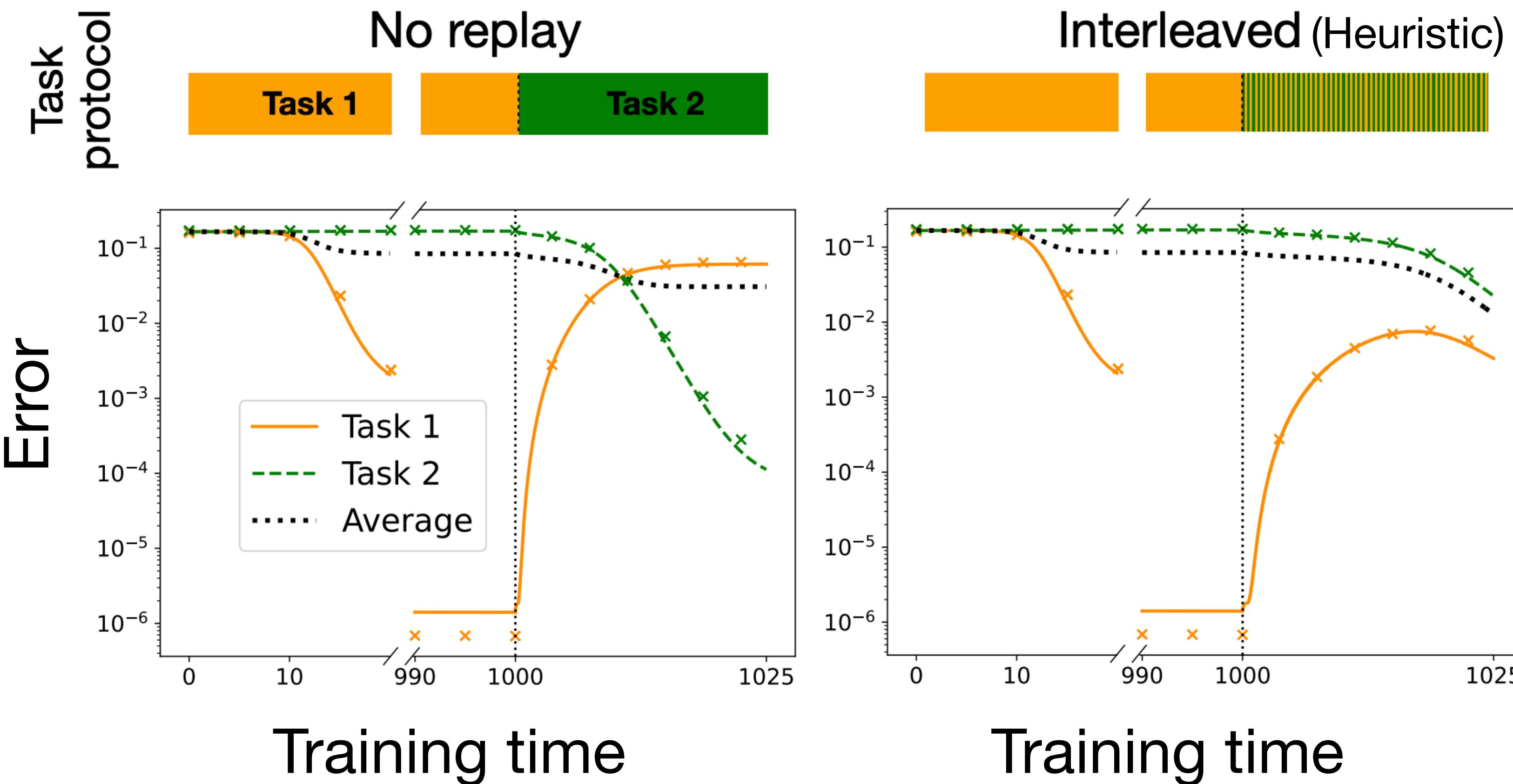
# Optimal strategy vs benchmarks

Control:  $\mathbf{u} = \text{task}$



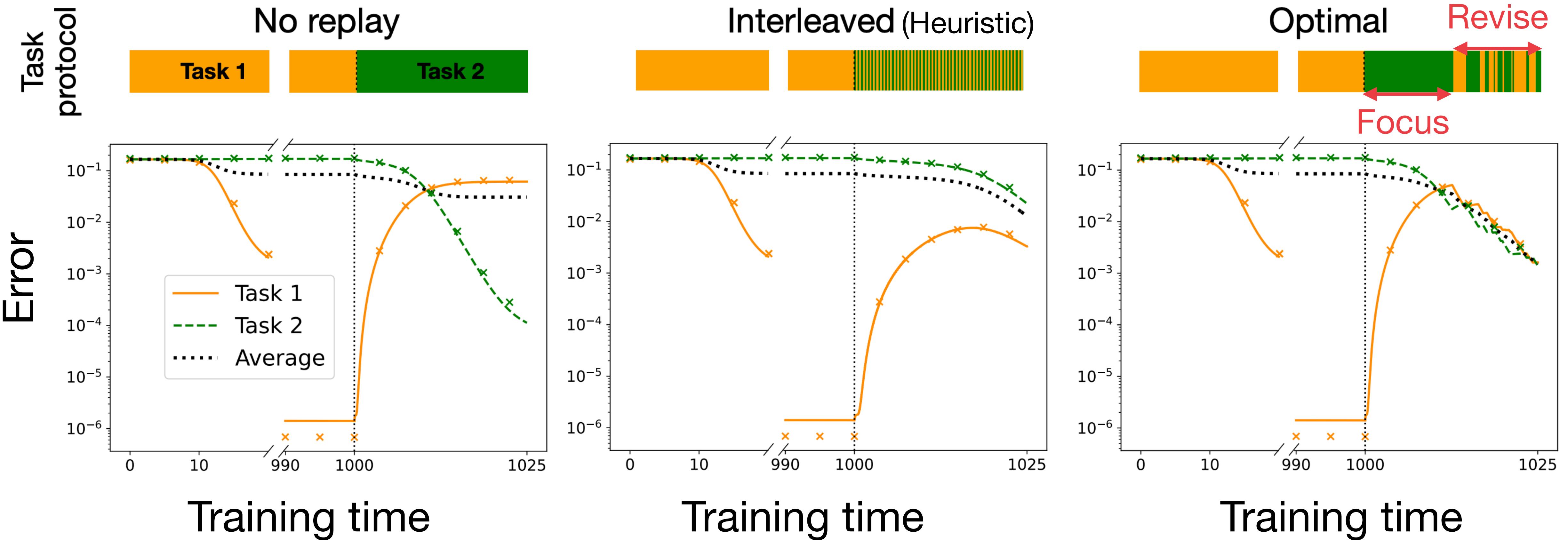
# Optimal strategy vs benchmarks

Control:  $\mathbf{u} = \text{task}$



# Optimal strategy vs benchmarks

Control:  $\mathbf{u} = \text{task}$

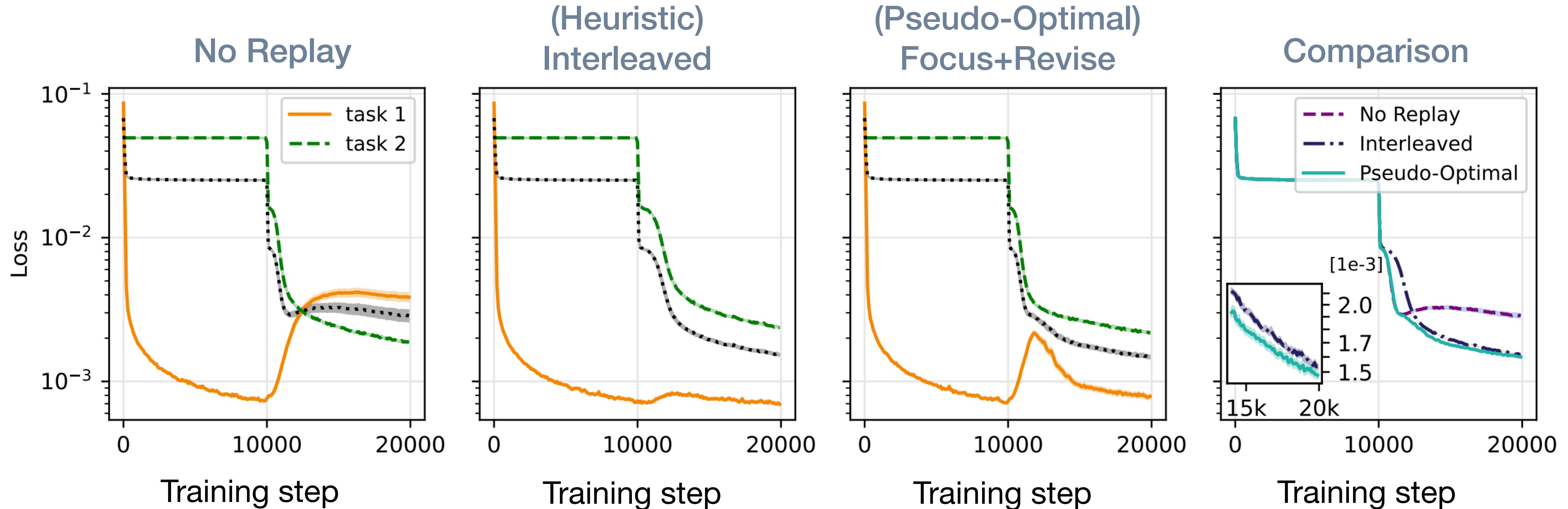


# Pseudo-optimal strategy on real data

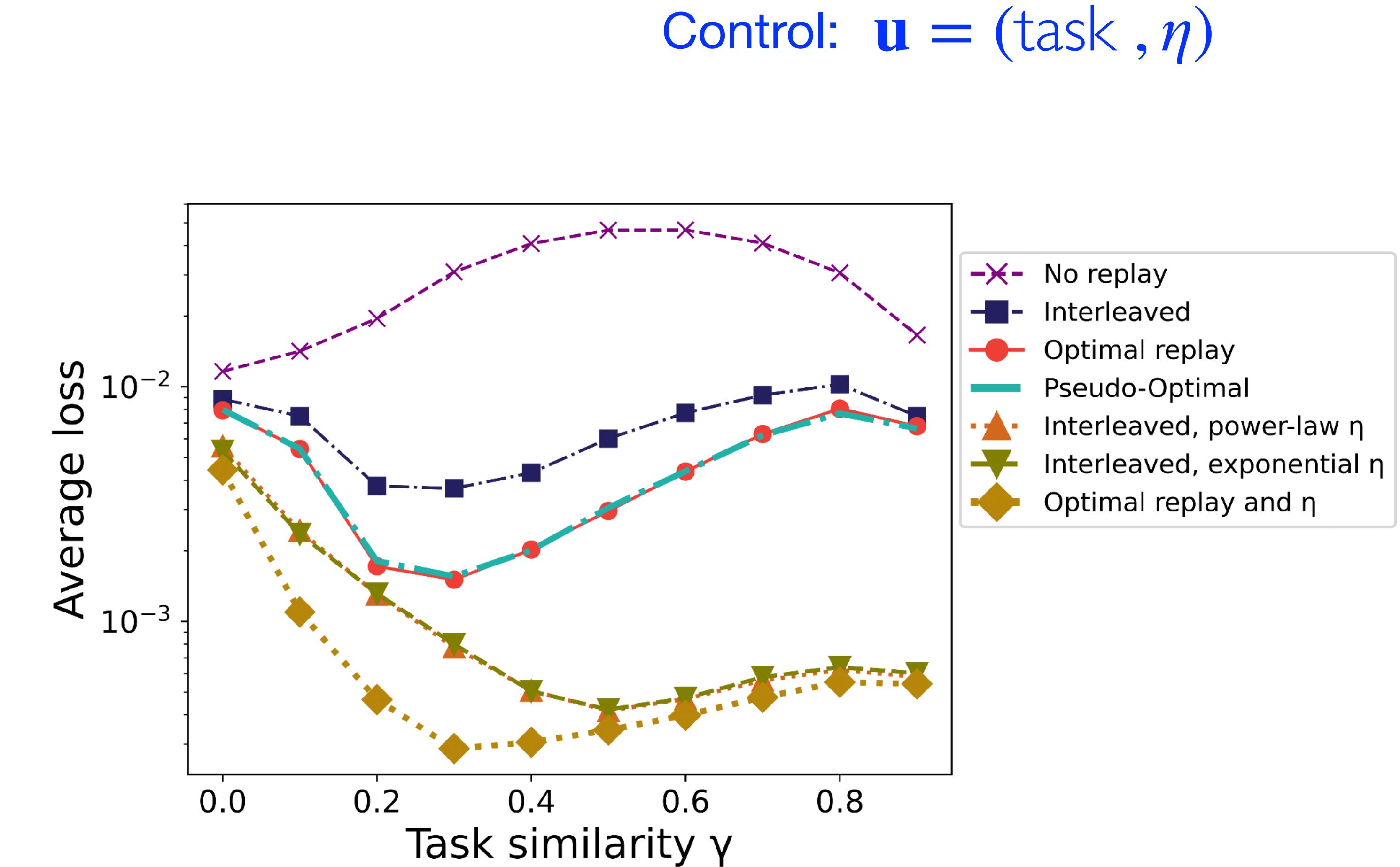
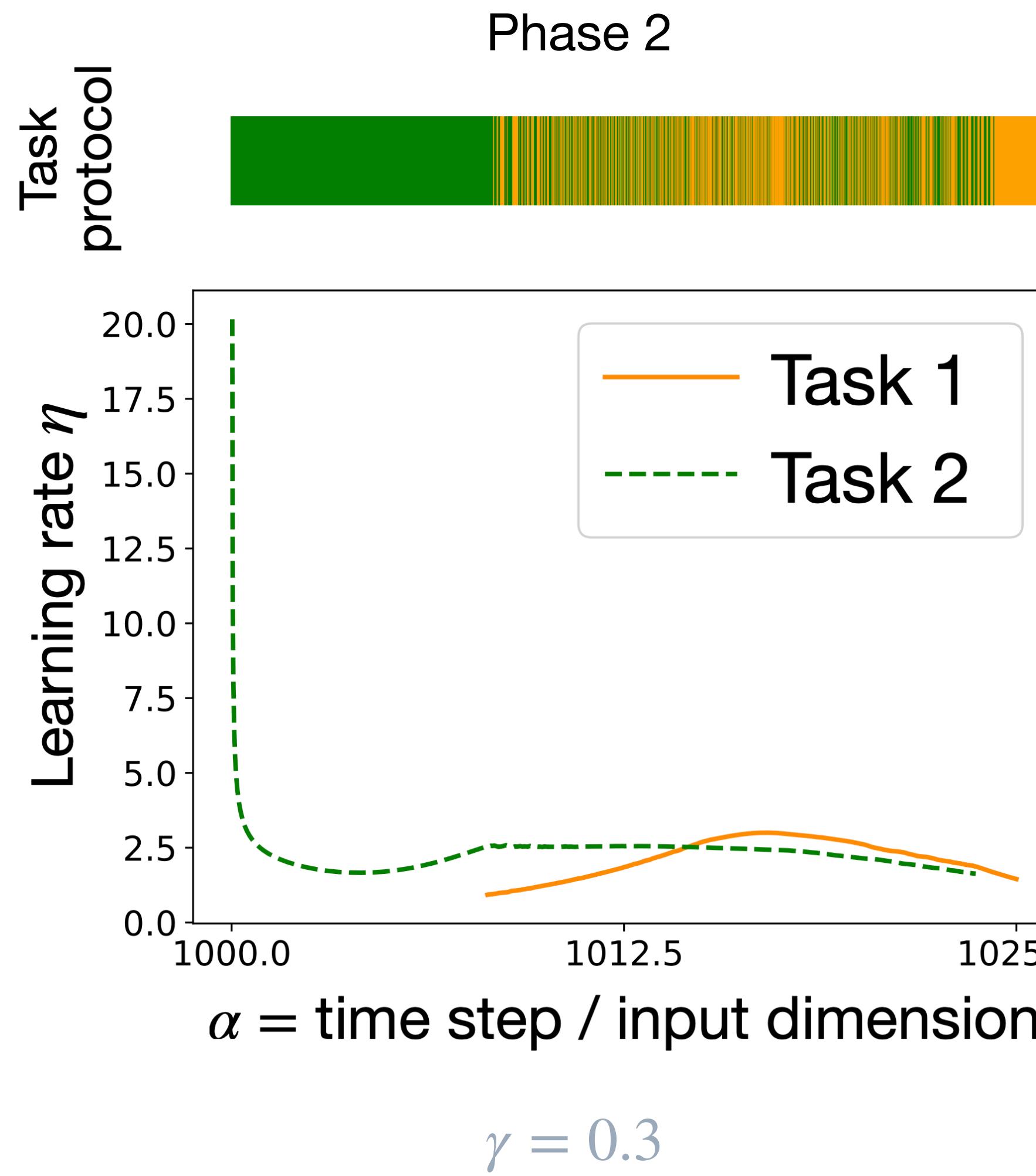
Task 1:  $\mathcal{D}_1 = \{\mathbf{x}_i^{(1)}, y_i^{(1)}\}_i$

Task 2:  $\mathcal{D}_2 = \{\mathbf{x}_i^{(2)}, y_i^{(2)}\}_i = \{\gamma \mathbf{x}_i^{(1)} + (1 - \gamma) \tilde{\mathbf{x}}_i, \gamma y_i^{(1)} + (1 - \gamma) \tilde{y}_i\}_i \quad \gamma = 0.5$

## Experiments on Fashion-MNIST:



# Joint optimization of replay and learning rate schedule



# Part III: *Denoising autoencoders*

# Denoising with autoencoders

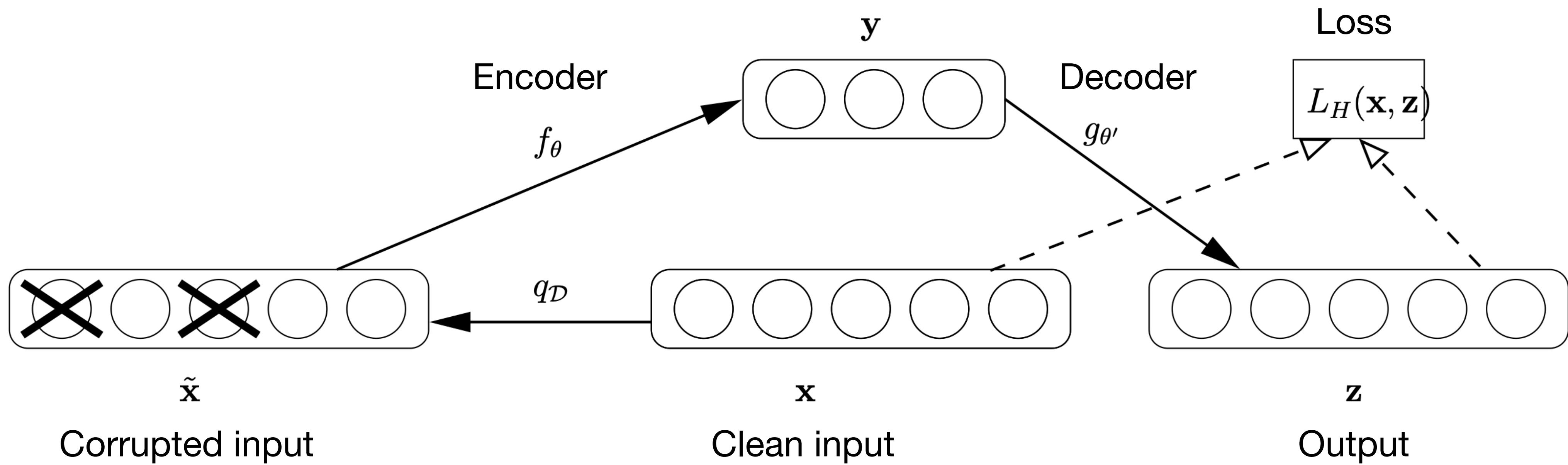


Image source: Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. *ICML 2008*

# Denoising with autoencoders

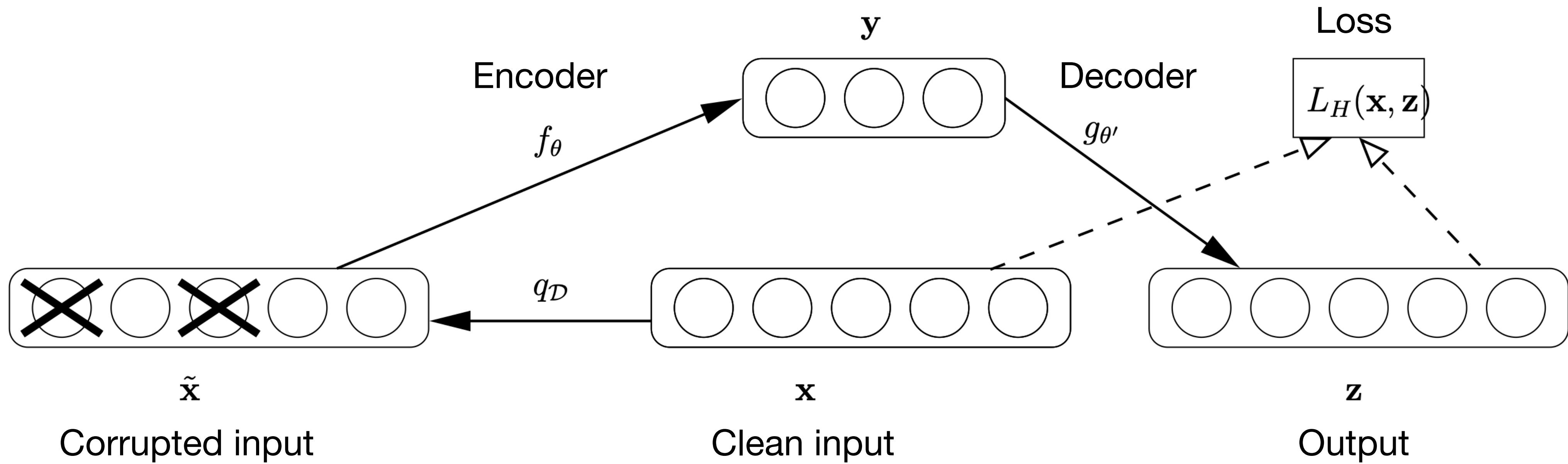


Image source: Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. ICML 2008

Theoretical works (non exhaustive list):

- Learning dynamics in linear DAEs [Pretorius, Kroon, and Kamper, ICML 2018]
- Exact asymptotics in nonlinear DAEs [Cui, and Zdeborová, NeurIPS 2023]
- Online learning in shallow reconstruction autoencoders [Refinetti, and Goldt, ICML 2022]
- Diffusion models parametrized by DAEs [Cui, Krzakala, Vanden-Eijnden, and Zdeborová, ICLR 2024; Cui, Pehlevan, and Lu, 2025]

# A prototypical model of DAE

See also: [Cui, and Zdeborová, NeurIPS 2023]

Clean input:  $\mathbf{x} \sim \sum_{c=1}^C p_c \mathcal{N}(\mu_c, \sigma_c^2 \mathbf{I}_N)$

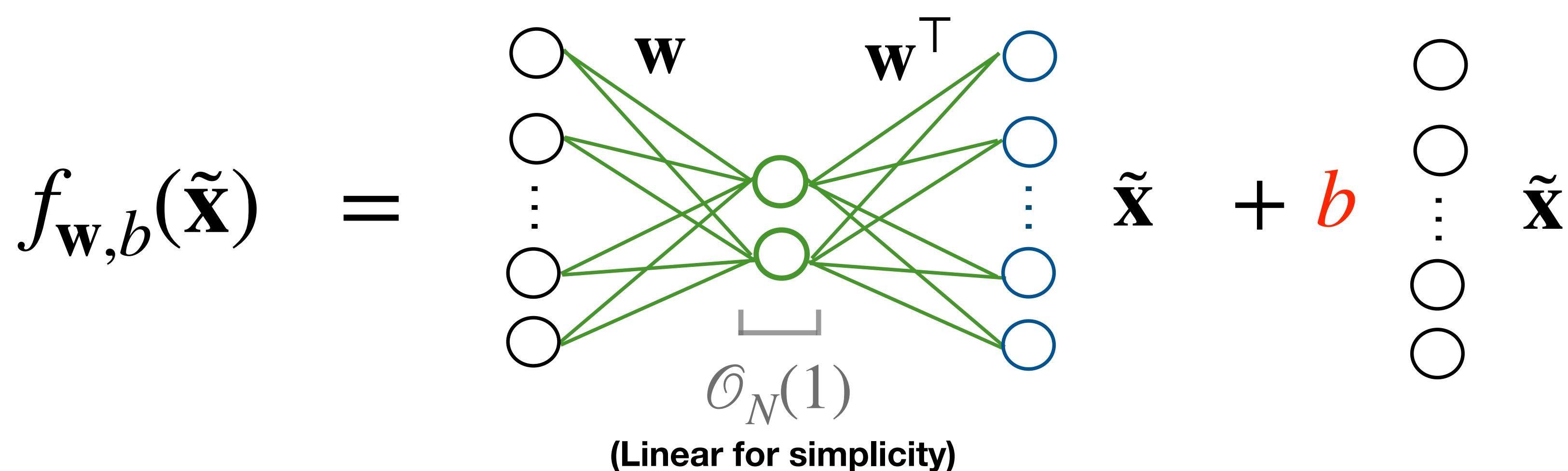
Noise level:  $\Delta \in (0,1)$

Corrupted input:  $\tilde{\mathbf{x}} = \sqrt{1 - \Delta} \mathbf{x} + \sqrt{\Delta} \xi$ ,  $\xi \sim \mathcal{N}(0, \mathbf{I}_N)$

**Two-layer DAE**

Bottleneck network

Skip connection

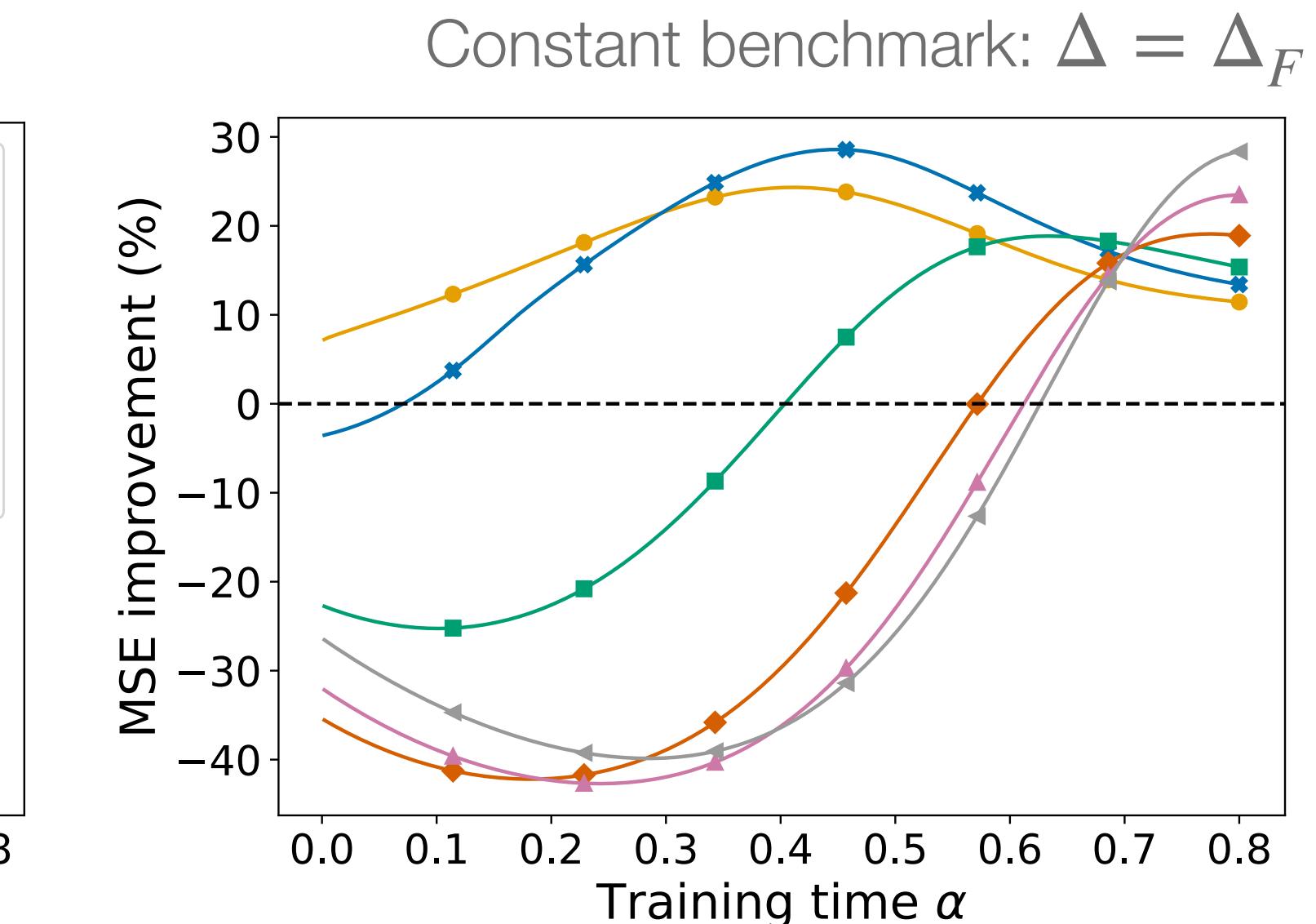
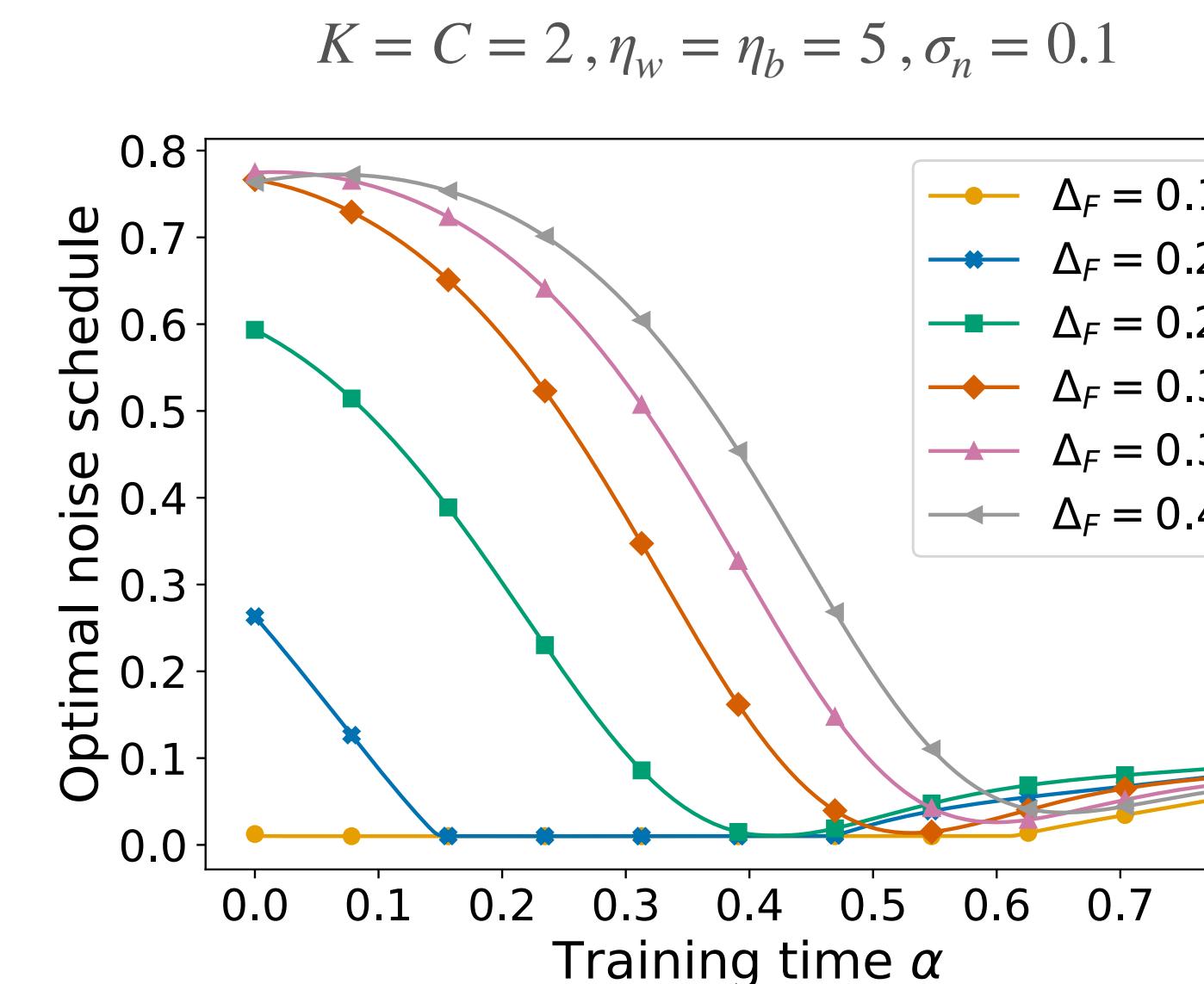


One-pass SGD on the loss:  $\mathcal{L}(\mathbf{w}, b) = \frac{1}{2} \| f_{\mathbf{w}}(\tilde{\mathbf{x}}) - \mathbf{x} \|^2$

# Optimal noise schedule at fixed test noise level

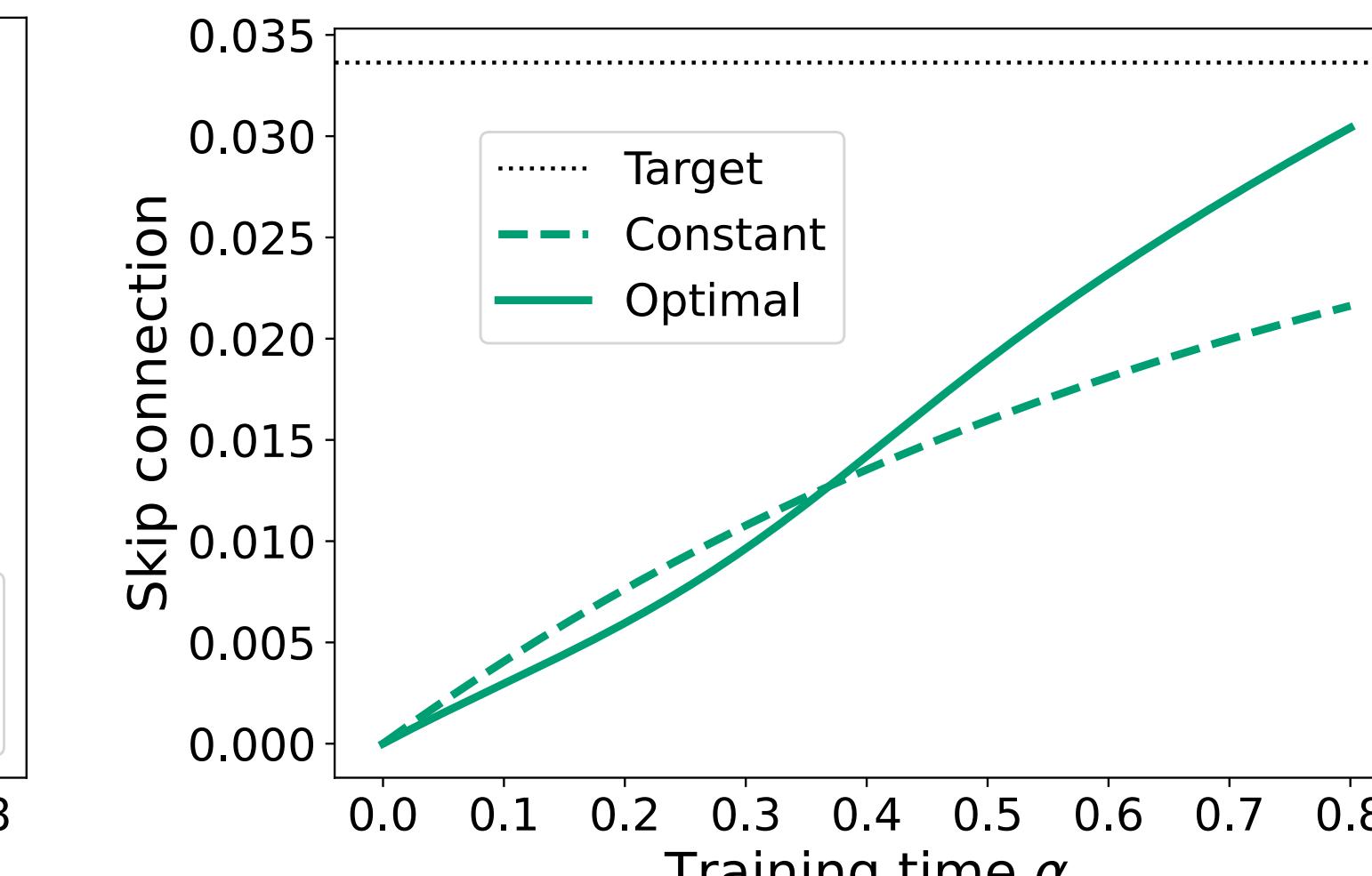
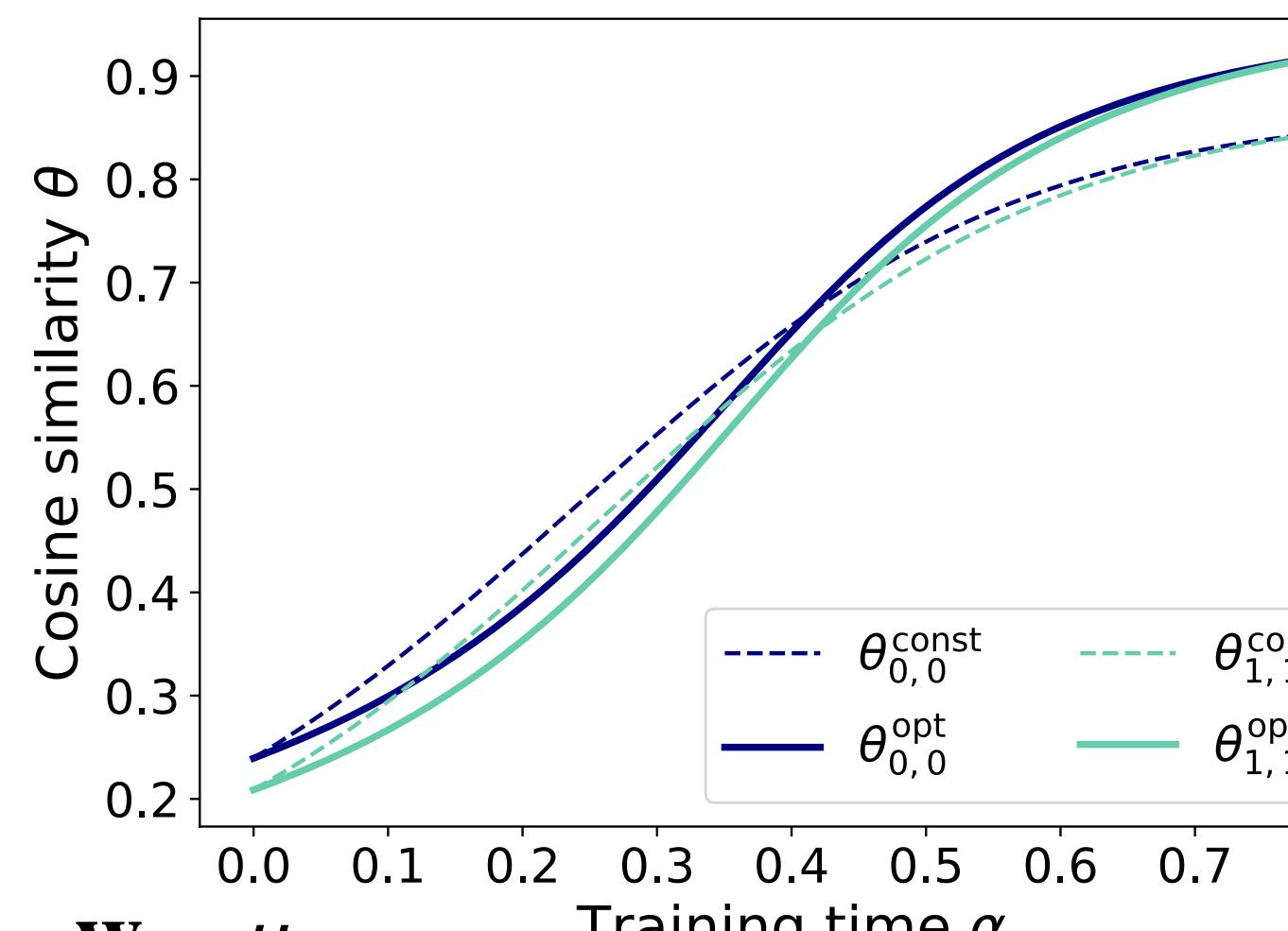
Control:  $u(\alpha) = \Delta(\alpha)$

Final performance: at  
fixed noise level  $\Delta_F$



The optimal schedule:

- Enhances the reconstruction capability of the bottleneck network
- Accelerates the convergence of the skip connection



$$\theta_{k,c} = \frac{\mathbf{w}_k \cdot \boldsymbol{\mu}_c}{\sqrt{\|\mathbf{w}_k\| \|\boldsymbol{\mu}_c\|}}$$

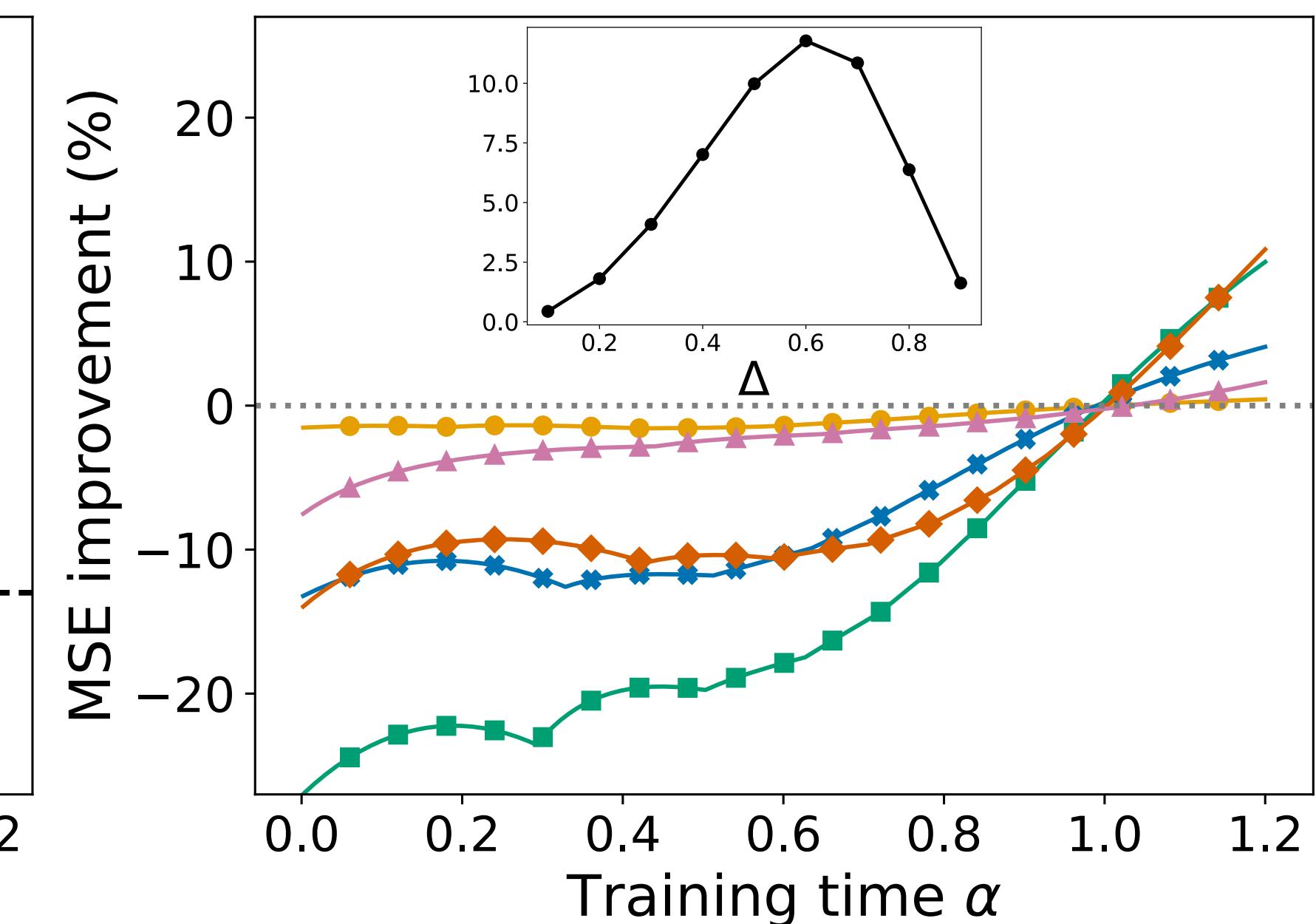
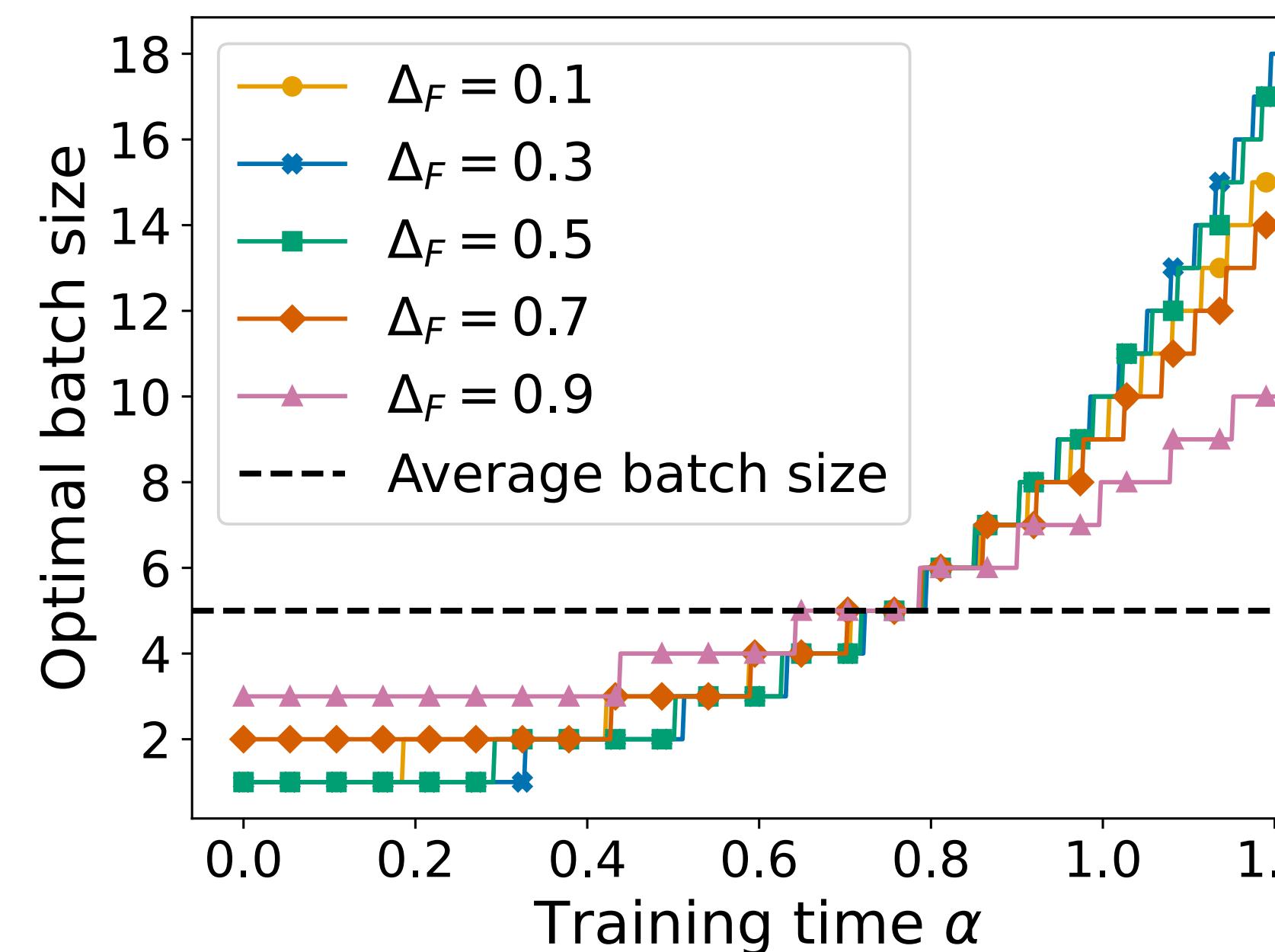
# Optimal batch size schedule

Batch augmentation:  $\tilde{\mathbf{x}}_a = \sqrt{1 - \Delta} \mathbf{x} + \sqrt{\Delta} \xi_a$ , with  $a = 1, \dots, B$  independent realizations of the noise

One-pass SGD on the loss:  $\mathcal{L}(\mathbf{w}) = \frac{1}{2B} \sum_{a=1}^B \|f_{\mathbf{w}}(\tilde{\mathbf{x}}_a) - \mathbf{x}\|^2$

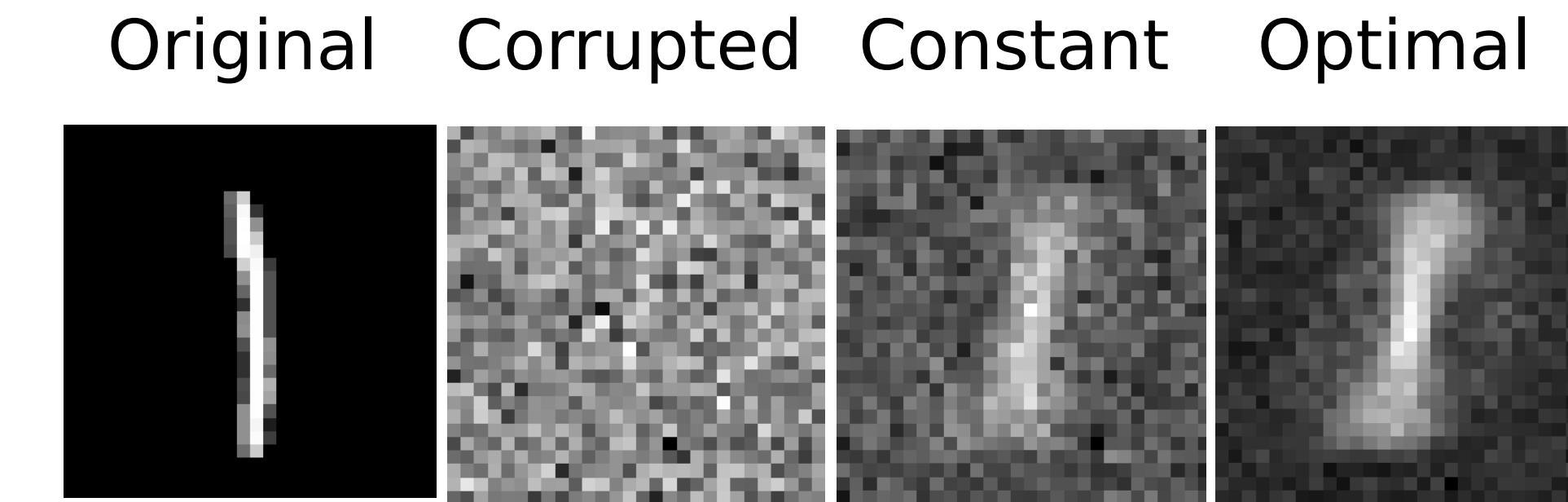
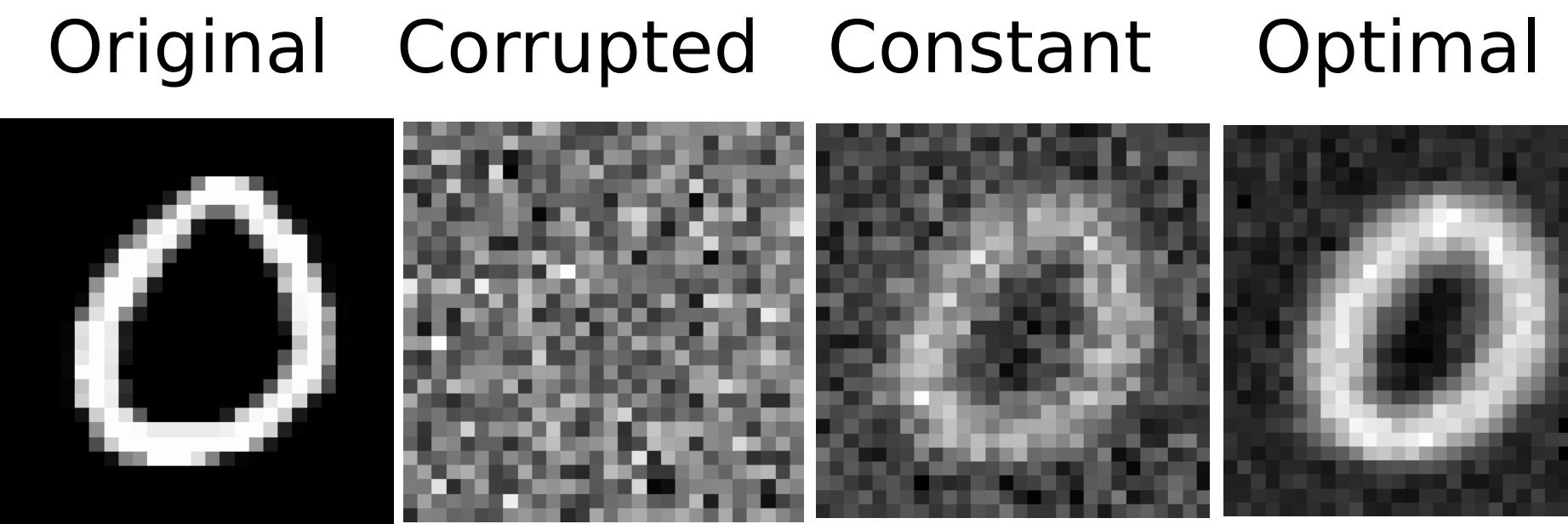
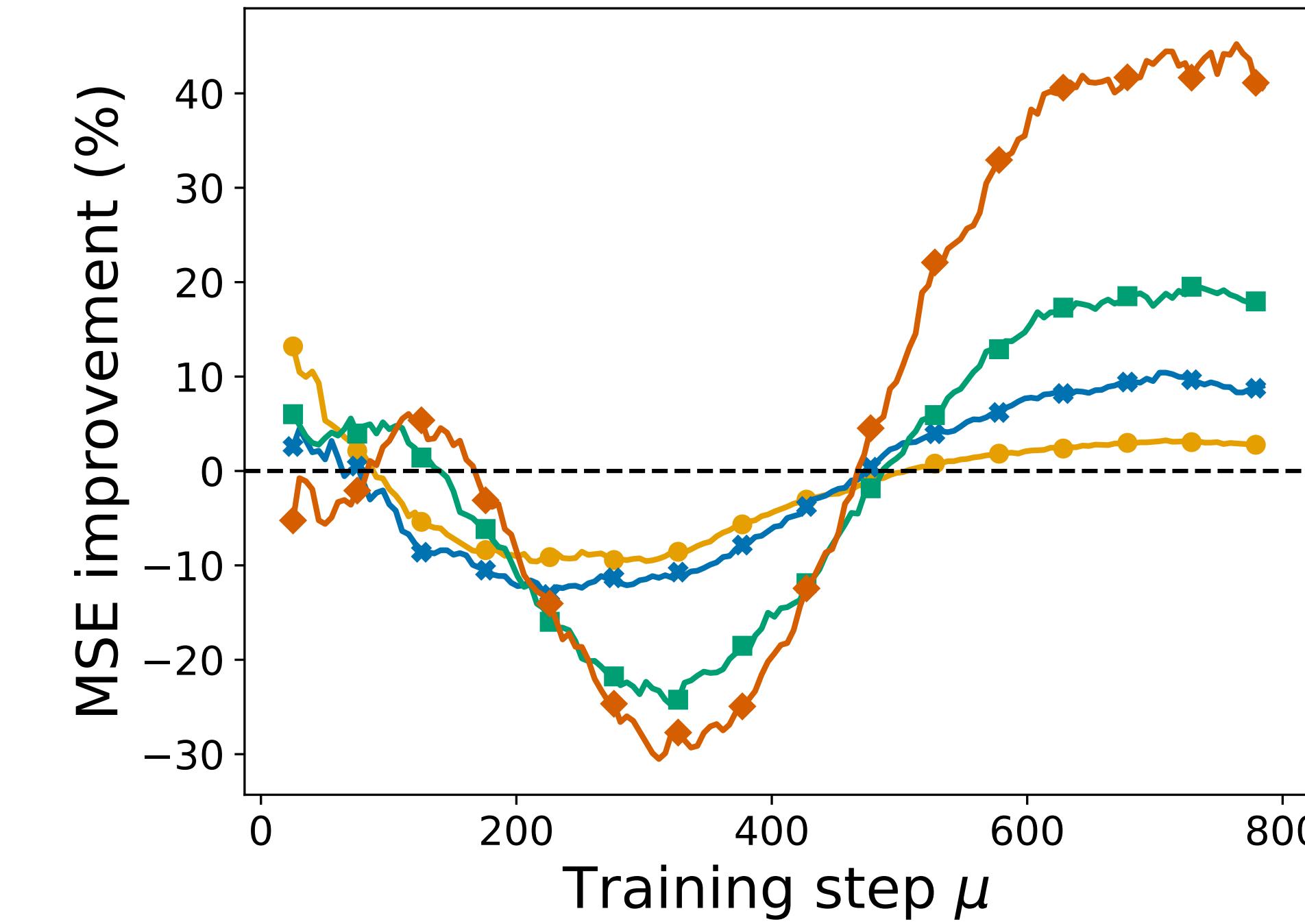
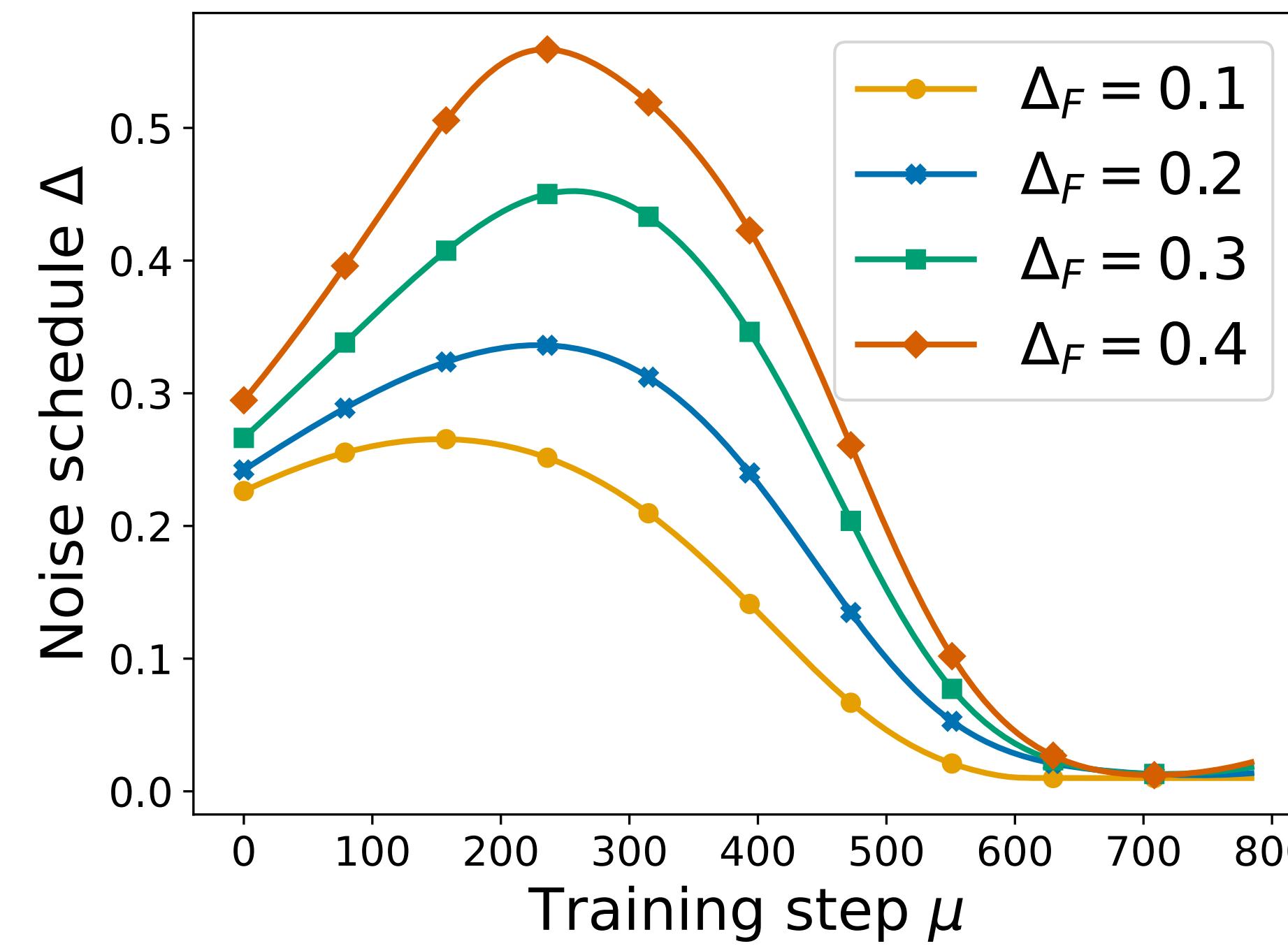
Constant benchmark:  $B = \bar{B}$

Control:  $u = B$   
 (time dependent  
 batch size, **fixed**  
**total budget**)



# MNIST experiments

Optimal noise schedules fitting only the cluster means and variances:



# Conclusions & Perspectives

- We formulate the design of learning protocols as an optimal control problem on the low-dimensional dynamics of the order parameters.
- Nontrivial yet interpretable strategies applicable to real data.
- Learning trade-offs

**Many open directions, e.g.,**

- Extend the theory to more realistic data models.
- Batch learning and memorization.
- Extend to different learning objectives (fairness metrics ...).

Thank you!