

Optimal learning strategies via statistical physics and control theory

Francesco Mori



ML Nosh Lunch

10th Feb 2025

“Optimal protocols for continual learning via
statistical physics and control theory”

FM, Stefano Sarao Mannelli, Francesca Mignacco

Accepted at ICLR 2025

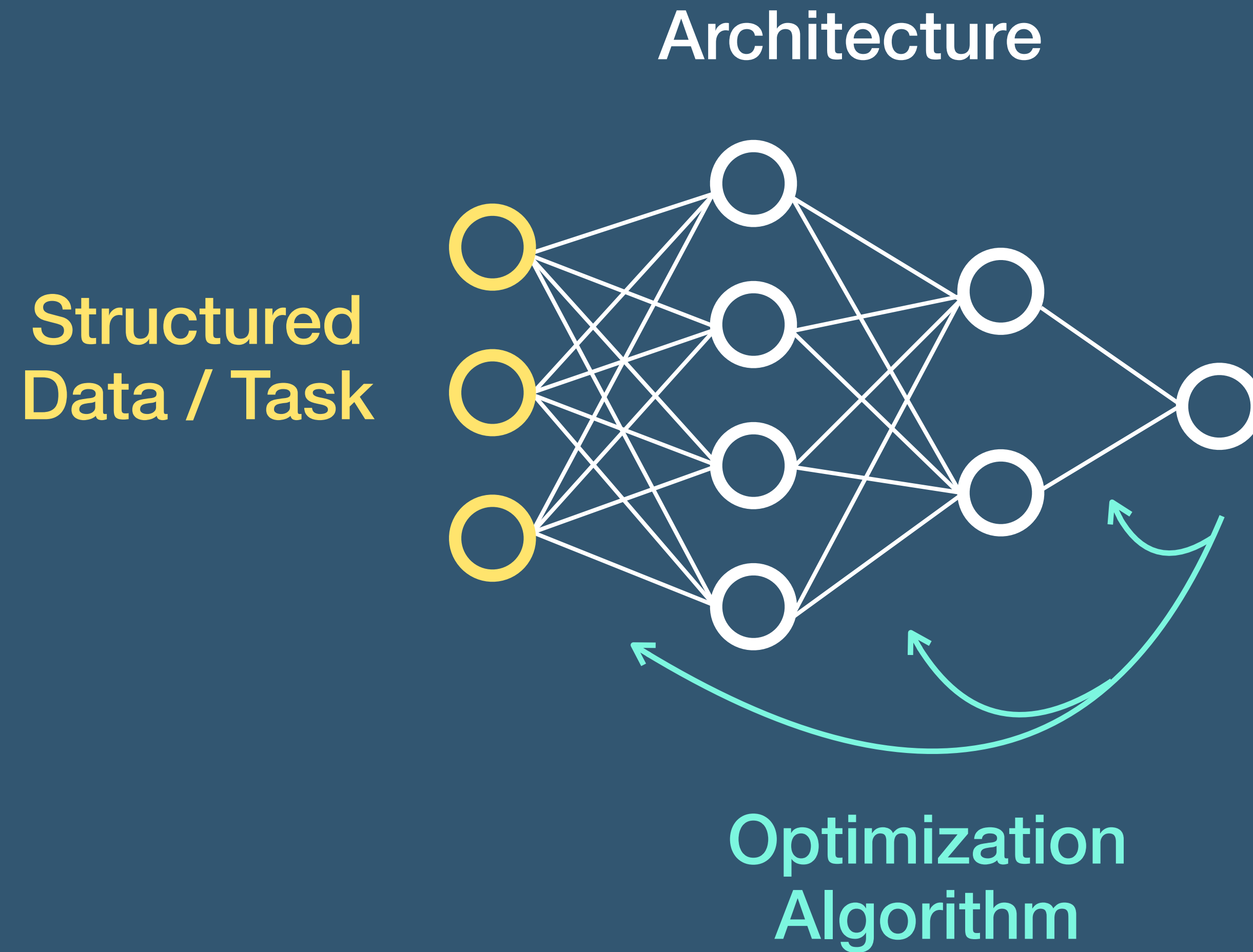


Francesca Mignacco
Princeton and CUNY

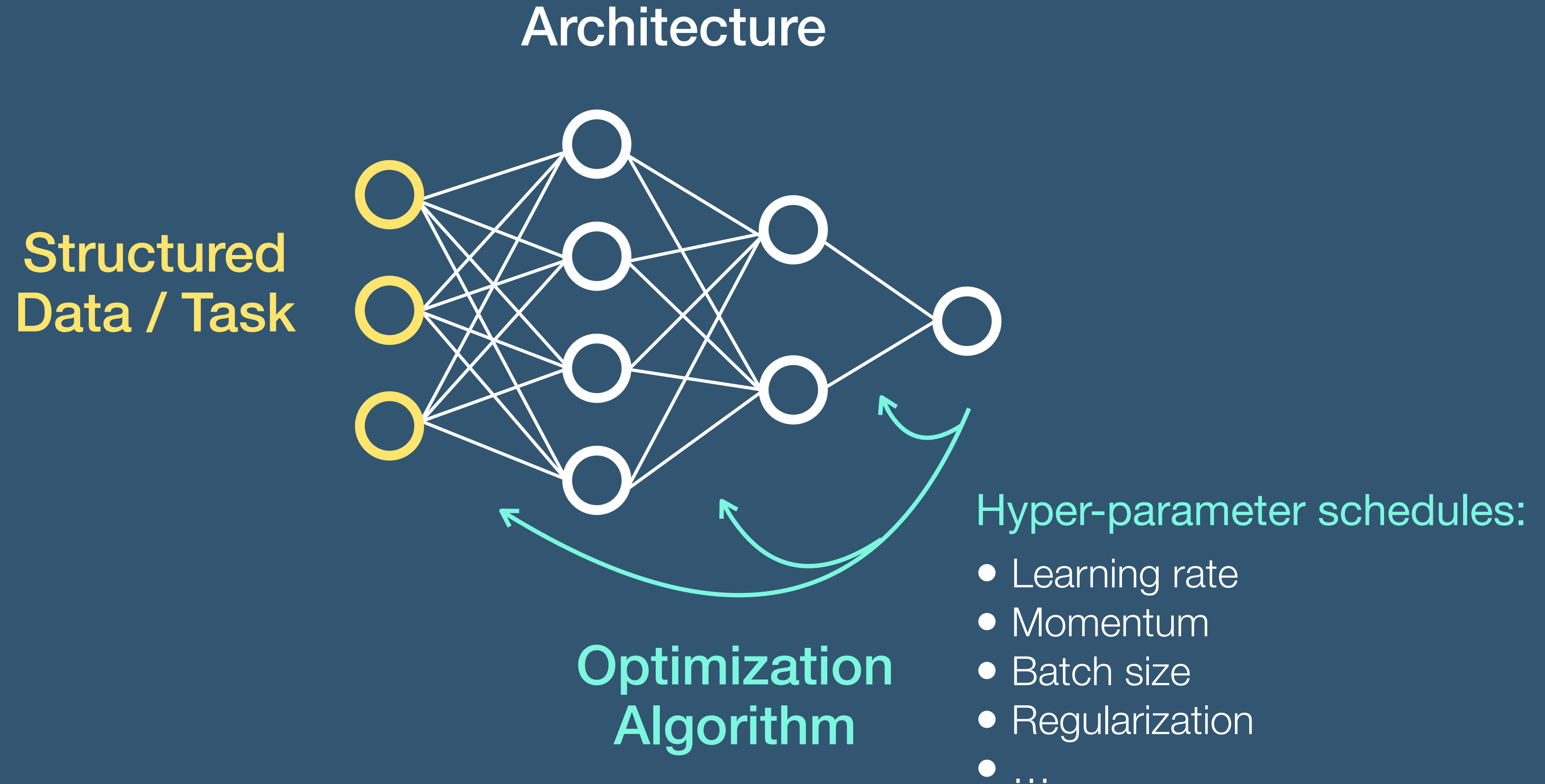


Stefano Sarao Mannelli
Chalmers University

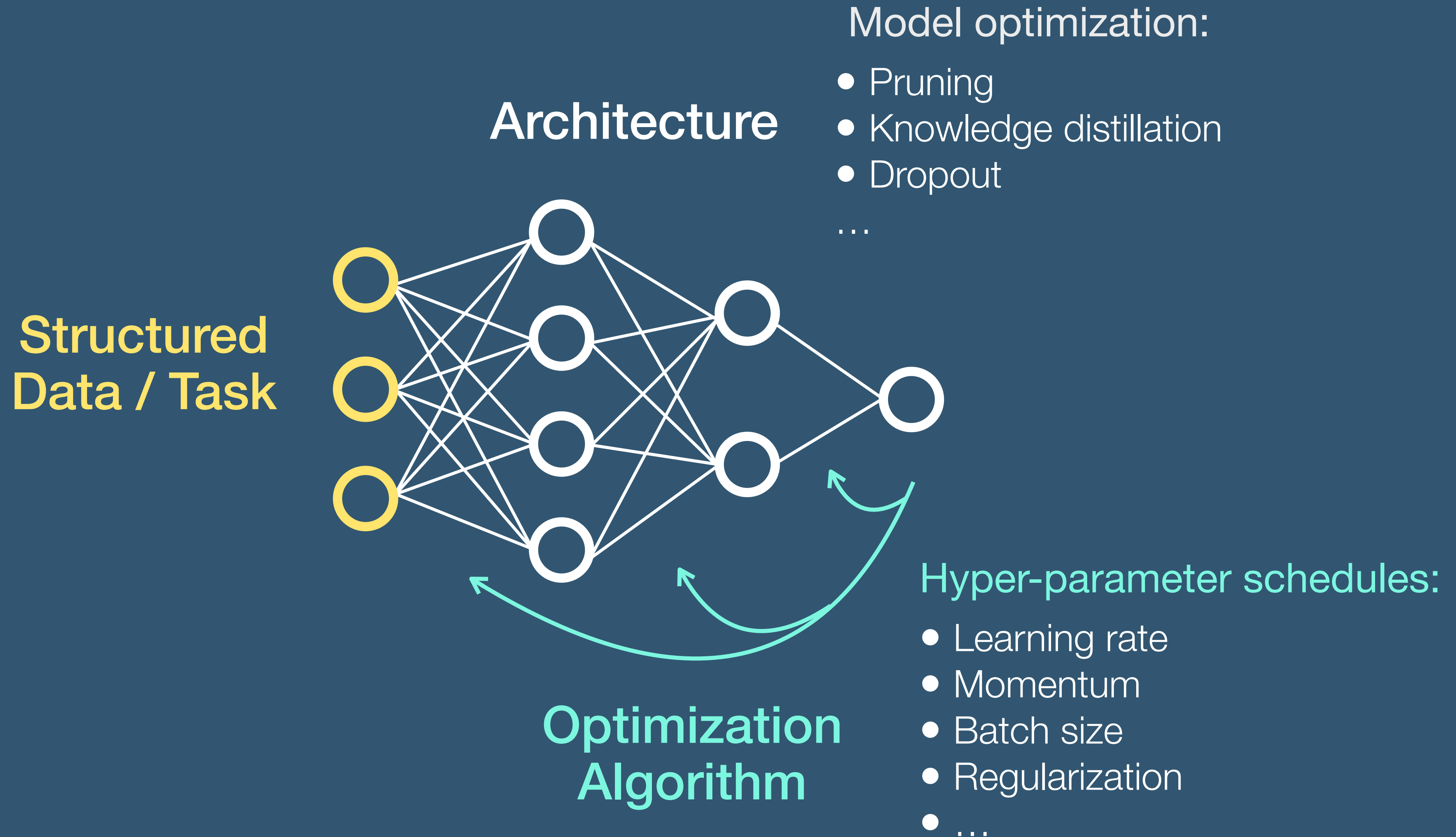
Learning protocols



Learning protocols

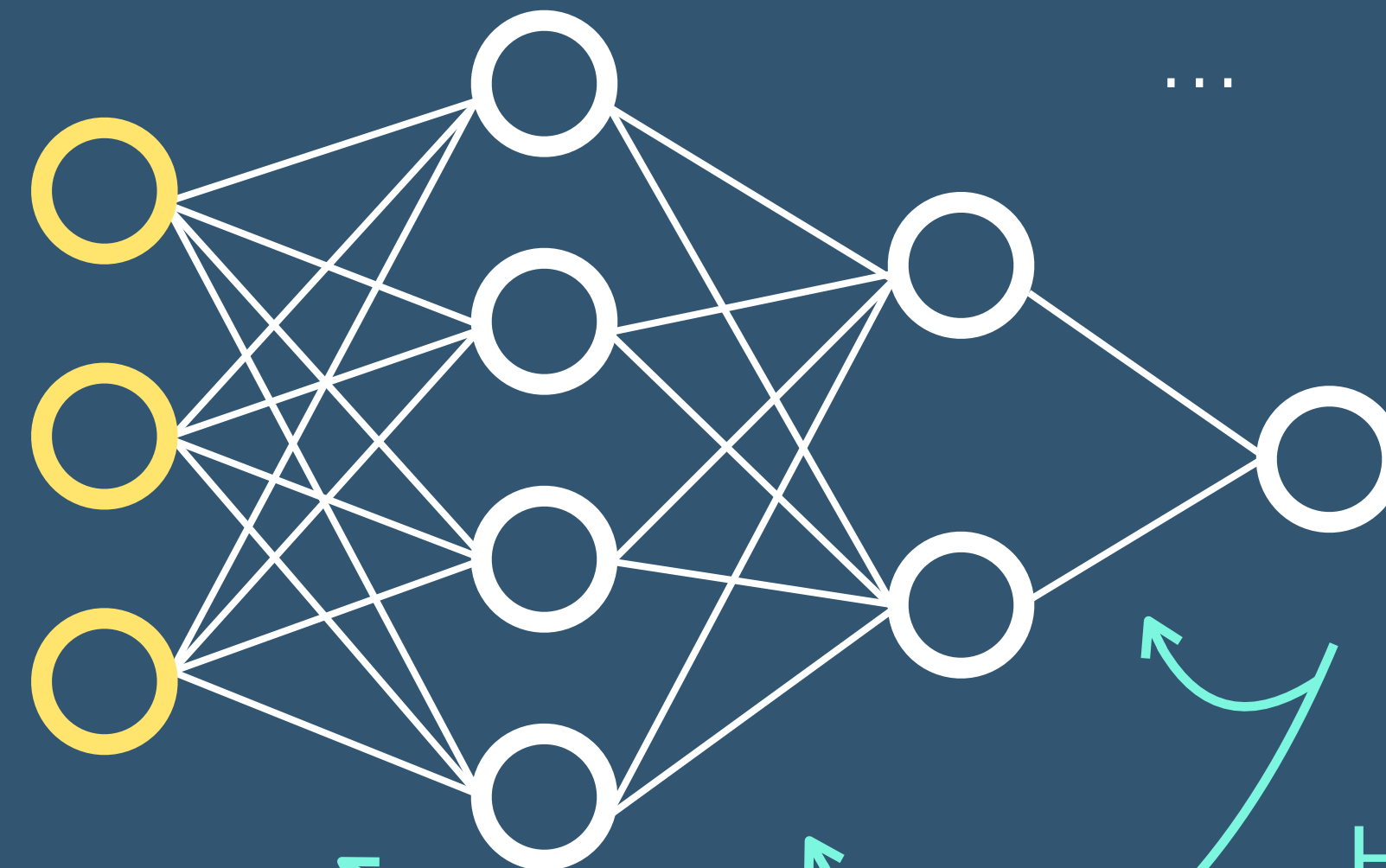


Learning protocols



Learning protocols

Structured
Data / Task



Model optimization:

- Pruning
- Knowledge distillation
- Dropout

...

Dynamic data / task selection:

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

Optimization
Algorithm

Hyper-parameter schedules:

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

Learning protocols

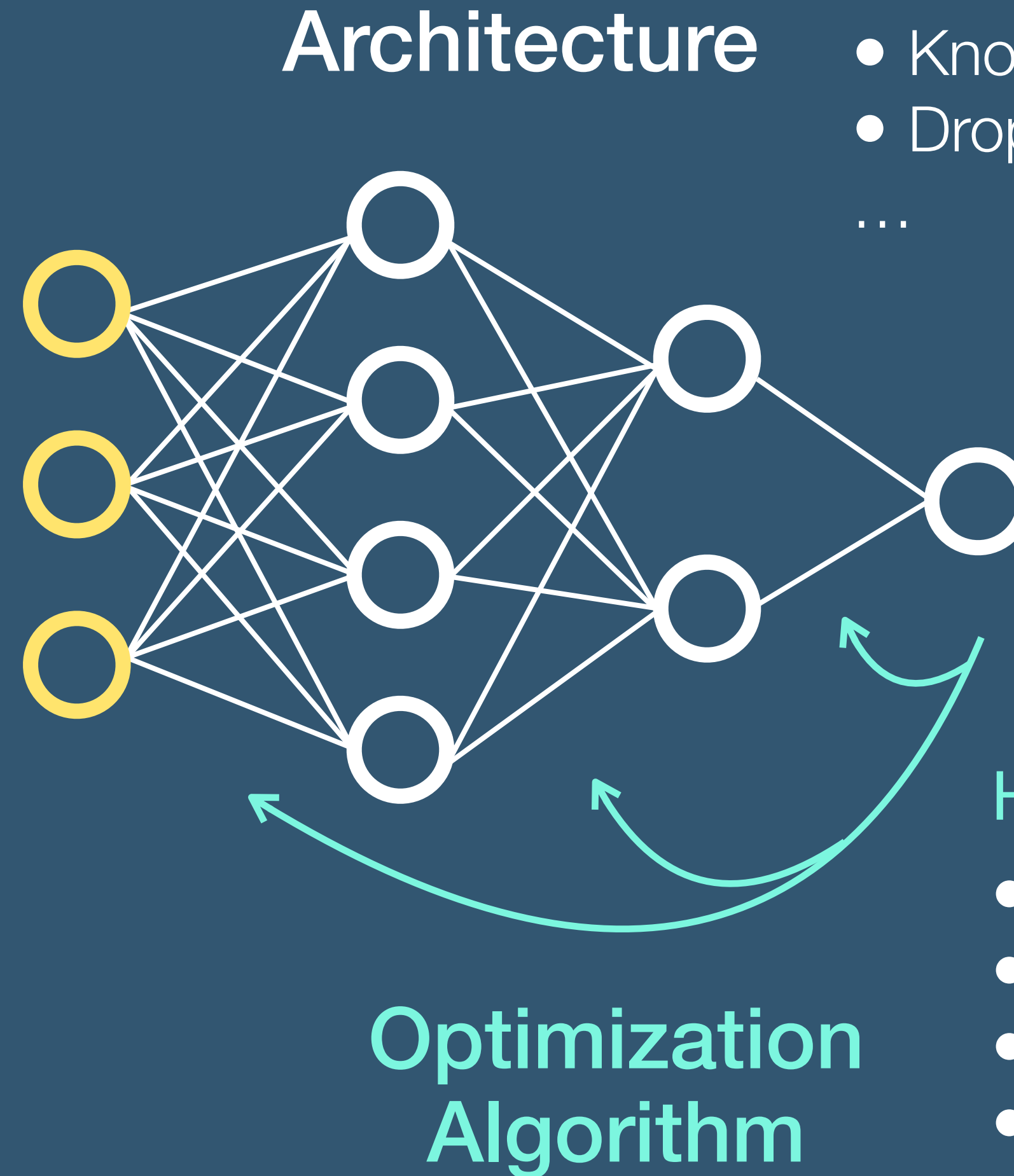
Goals:

- **Speed-up** convergence
- **Guide** the training towards better regions of parameters space

Model optimization:

- Pruning
- Knowledge distillation
- Dropout
- ...

Structured Data / Task



Dynamic data / task selection:

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

Hyper-parameter schedules:

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

Learning protocols

Goals:

- **Speed-up** convergence
- **Guide** the training towards better regions of parameters space

From smoother landscape to the target

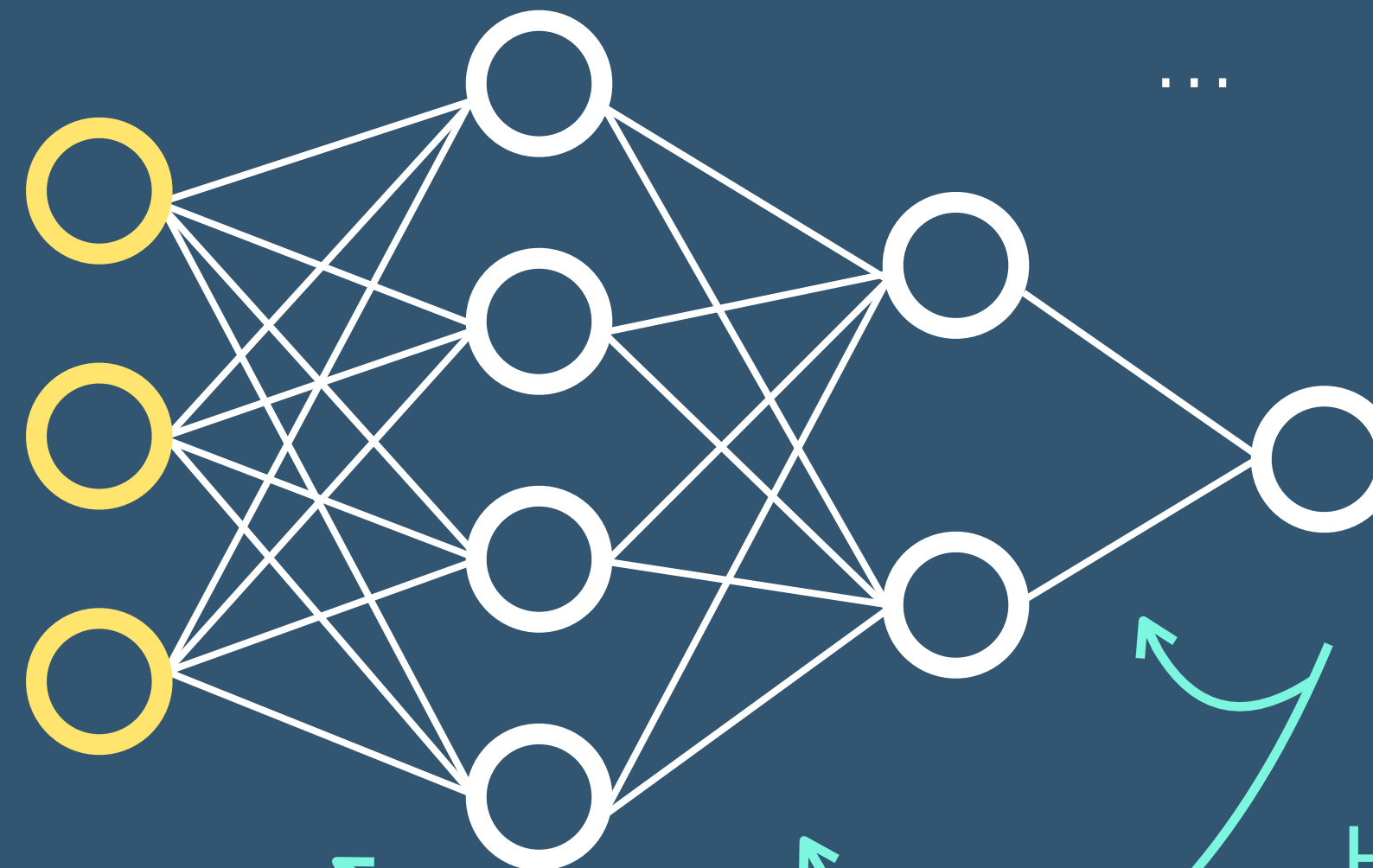
Structured Data / Task

In this talk:

Dynamic data / task selection:

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

Architecture



Optimization Algorithm

Model optimization:

- Pruning
- Knowledge distillation
- Dropout
- ...

Hyper-parameter schedules:

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

Dataset with labels

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^P$$

Neural Network

$$\hat{y} = f_{\mathbf{w}}(x)$$

Simplest example: $\hat{y} = \text{erf}(\mathbf{w}^T x)$

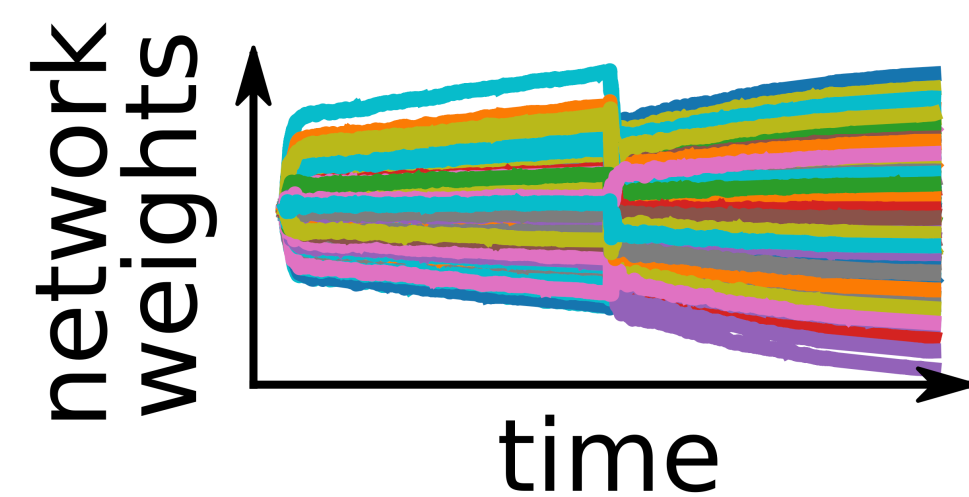
Error (aka loss)

$$\mathcal{L} = \frac{1}{2} (\hat{y} - y)^2$$

(Online) Stochastic gradient descent

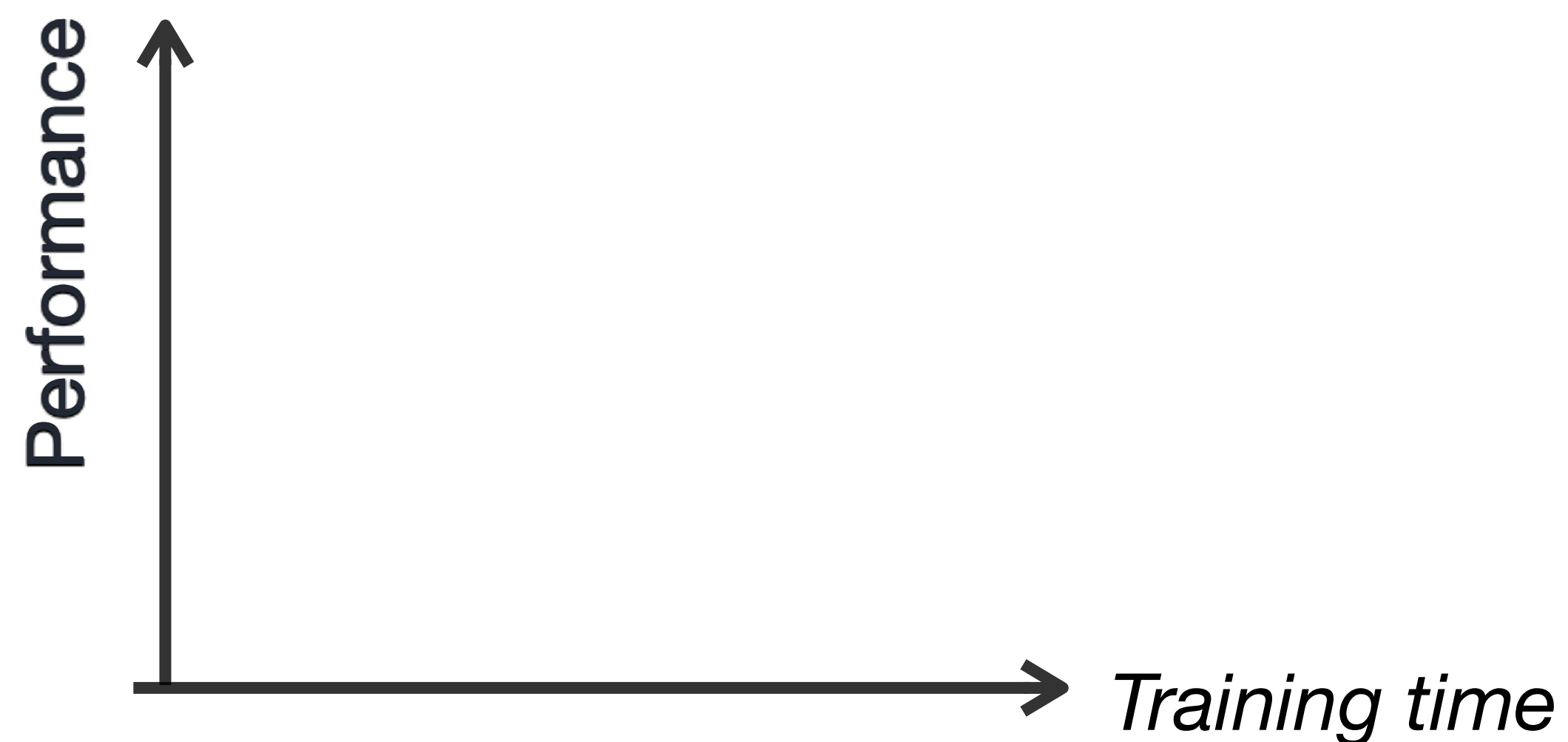
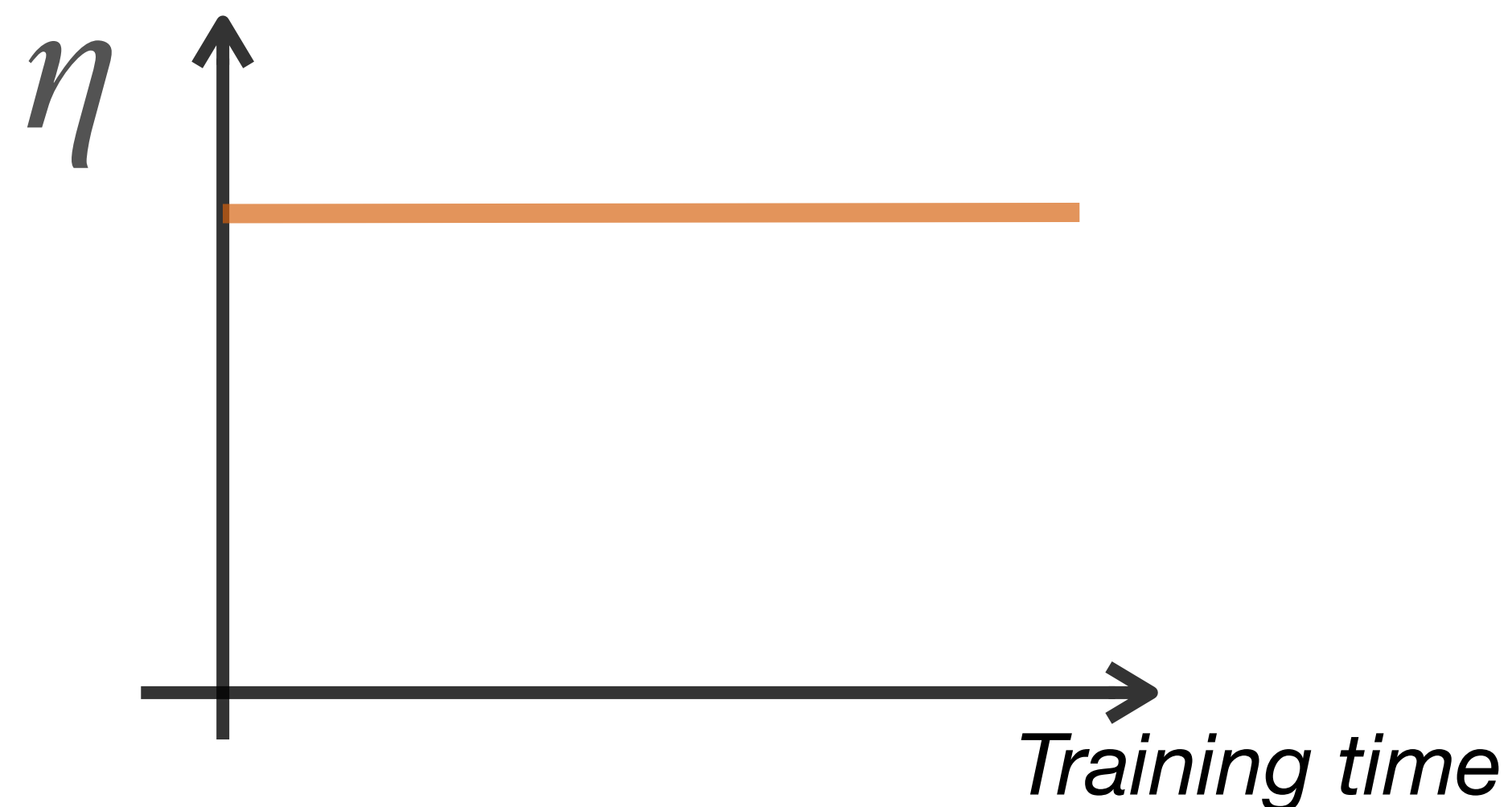
$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

Training protocols example: learning rate

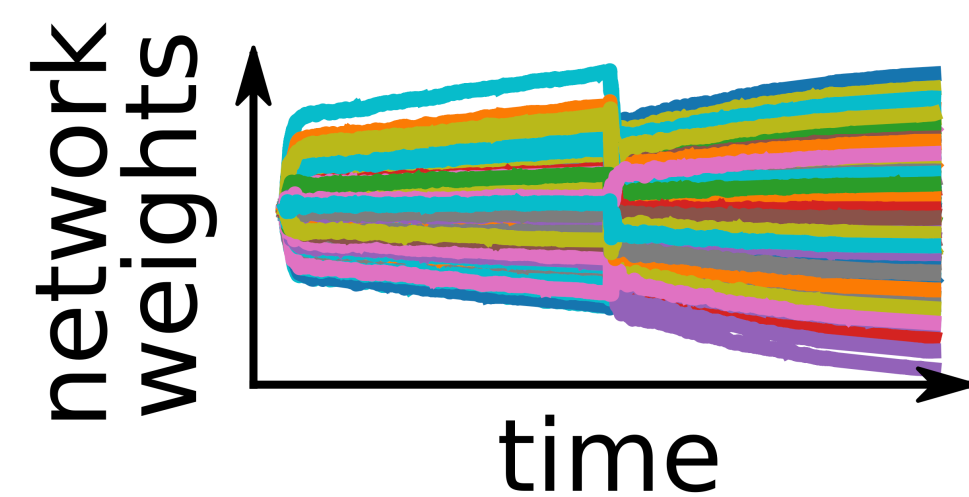


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

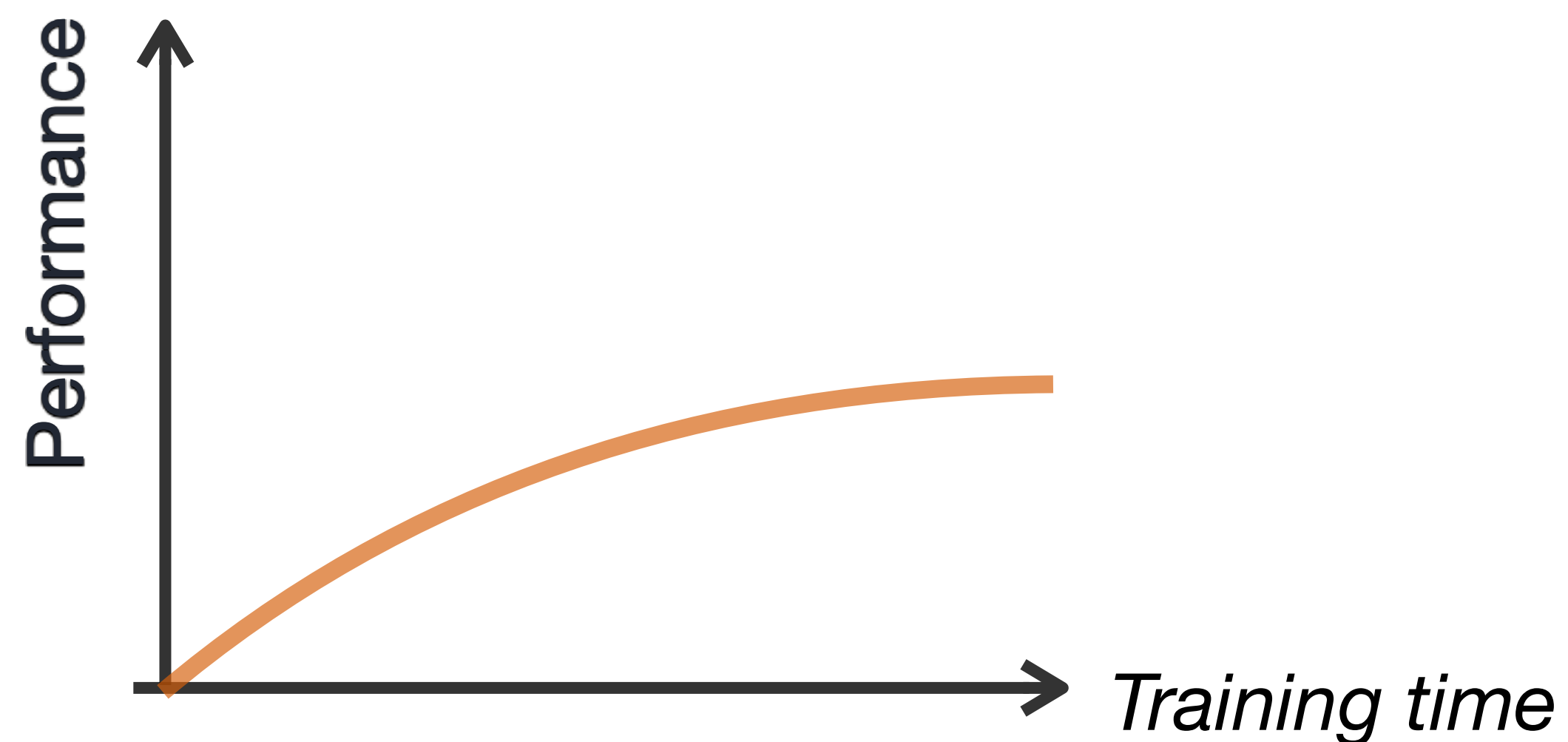
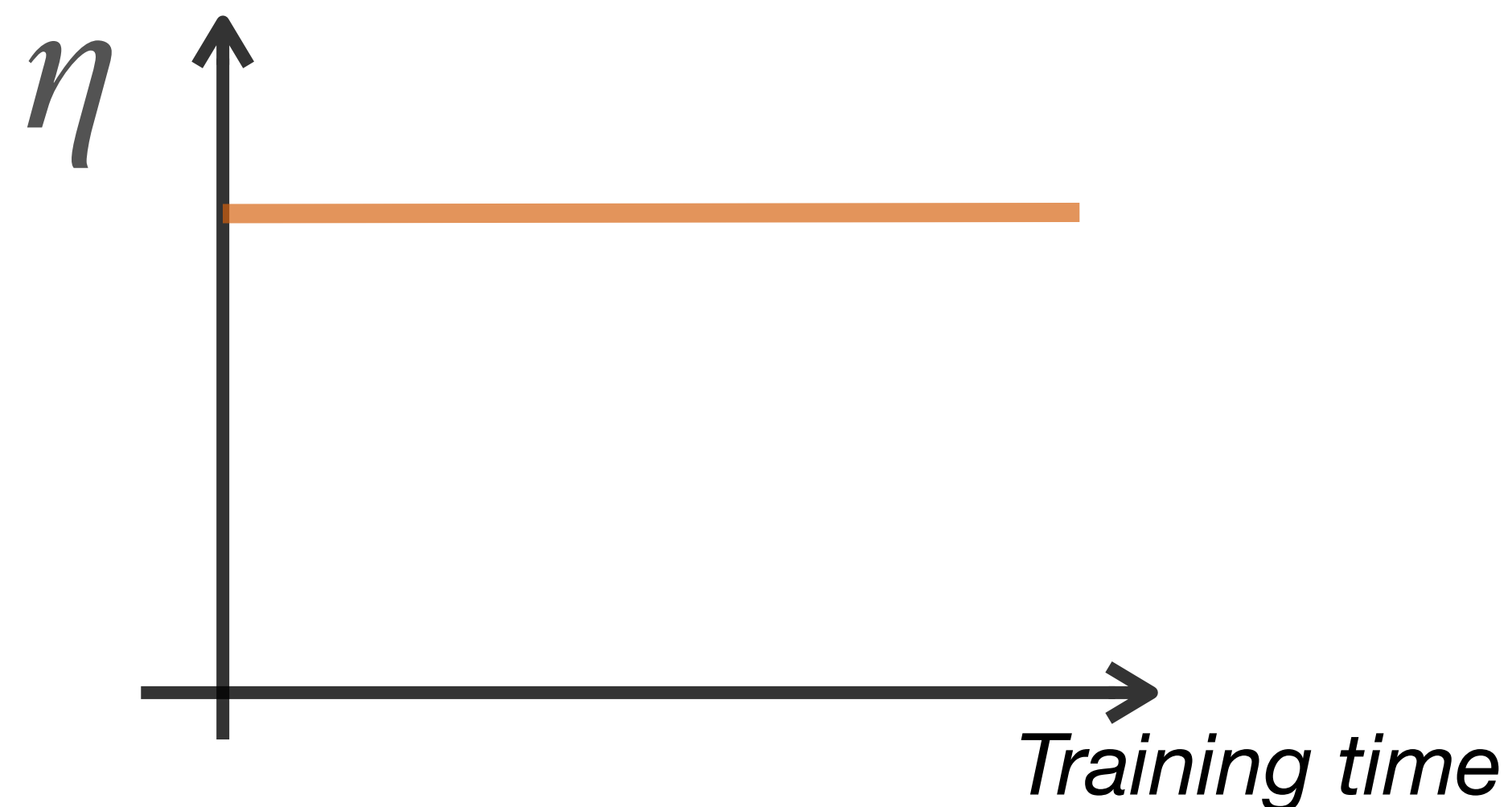


Training protocols example: learning rate

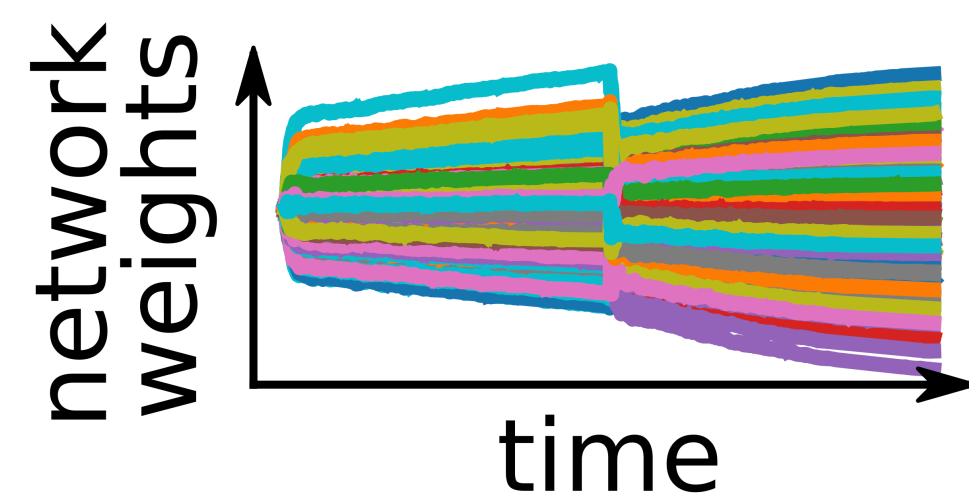


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

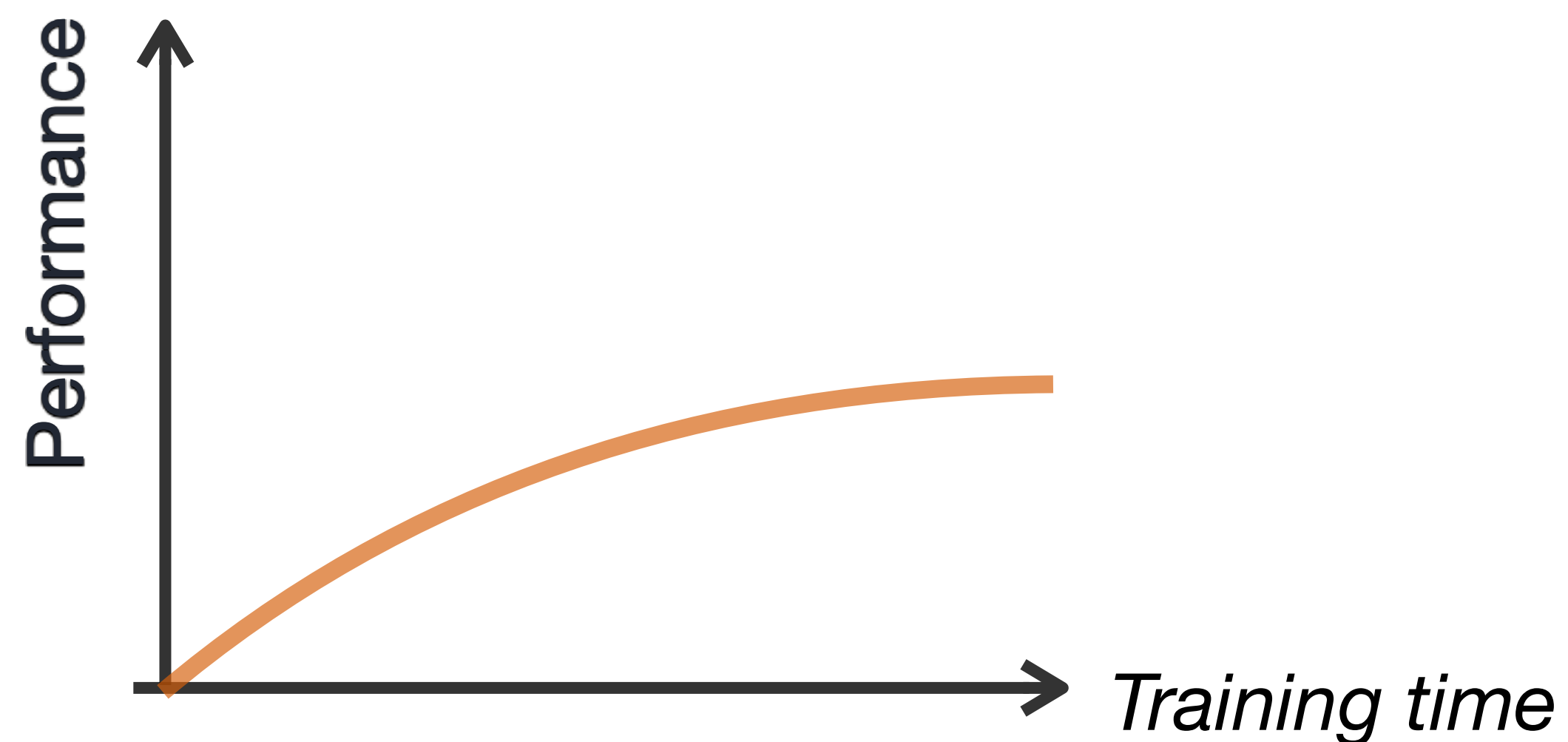
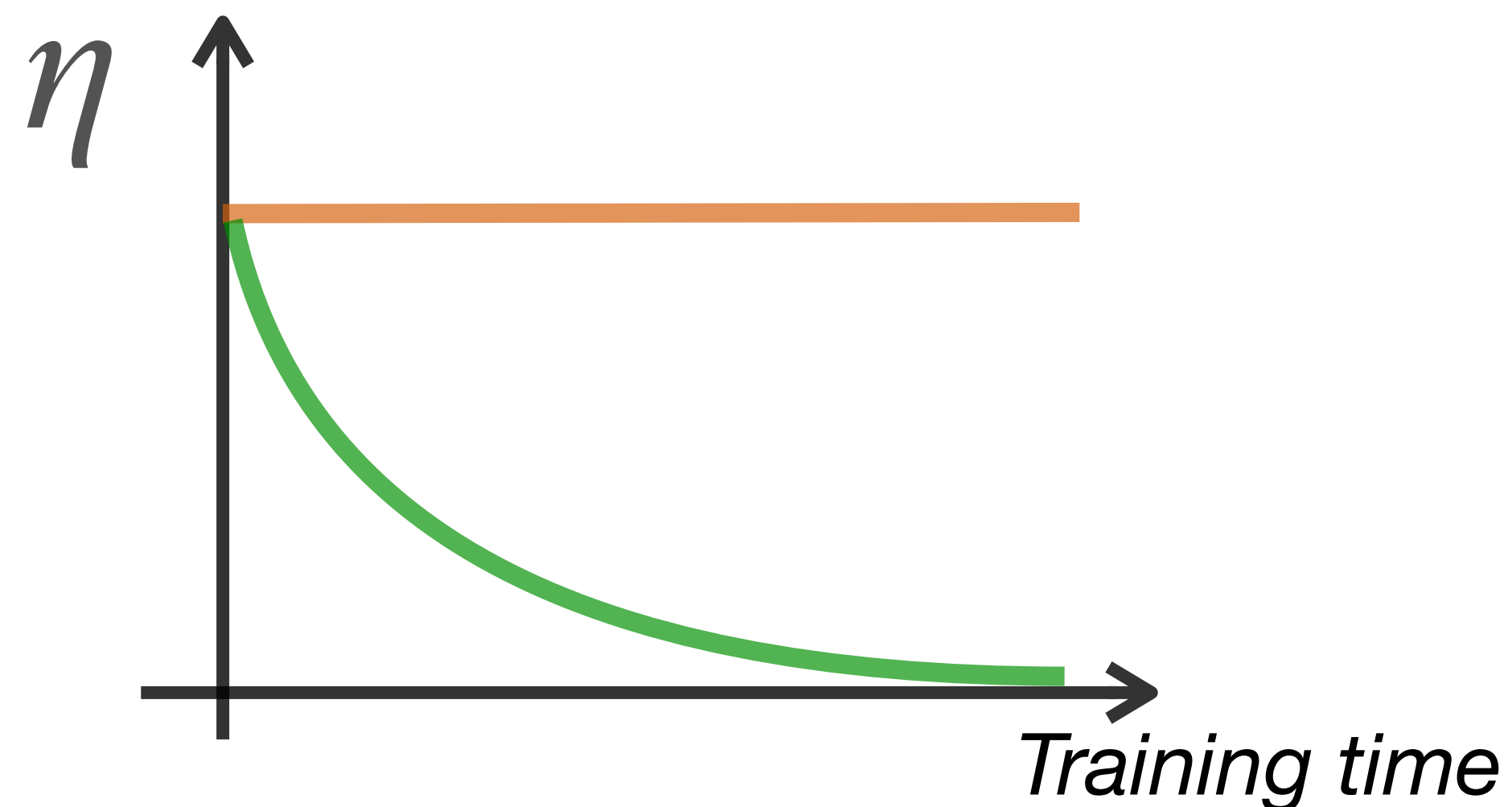


Training protocols example: learning rate

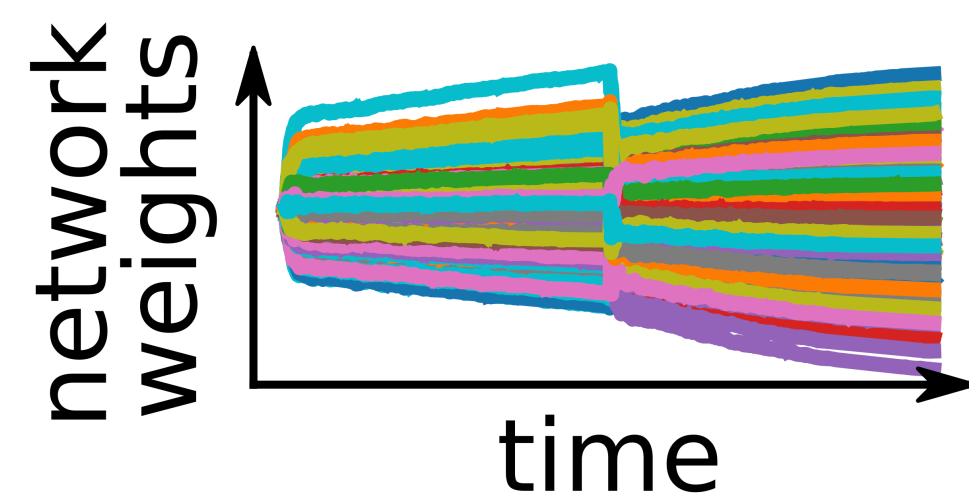


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

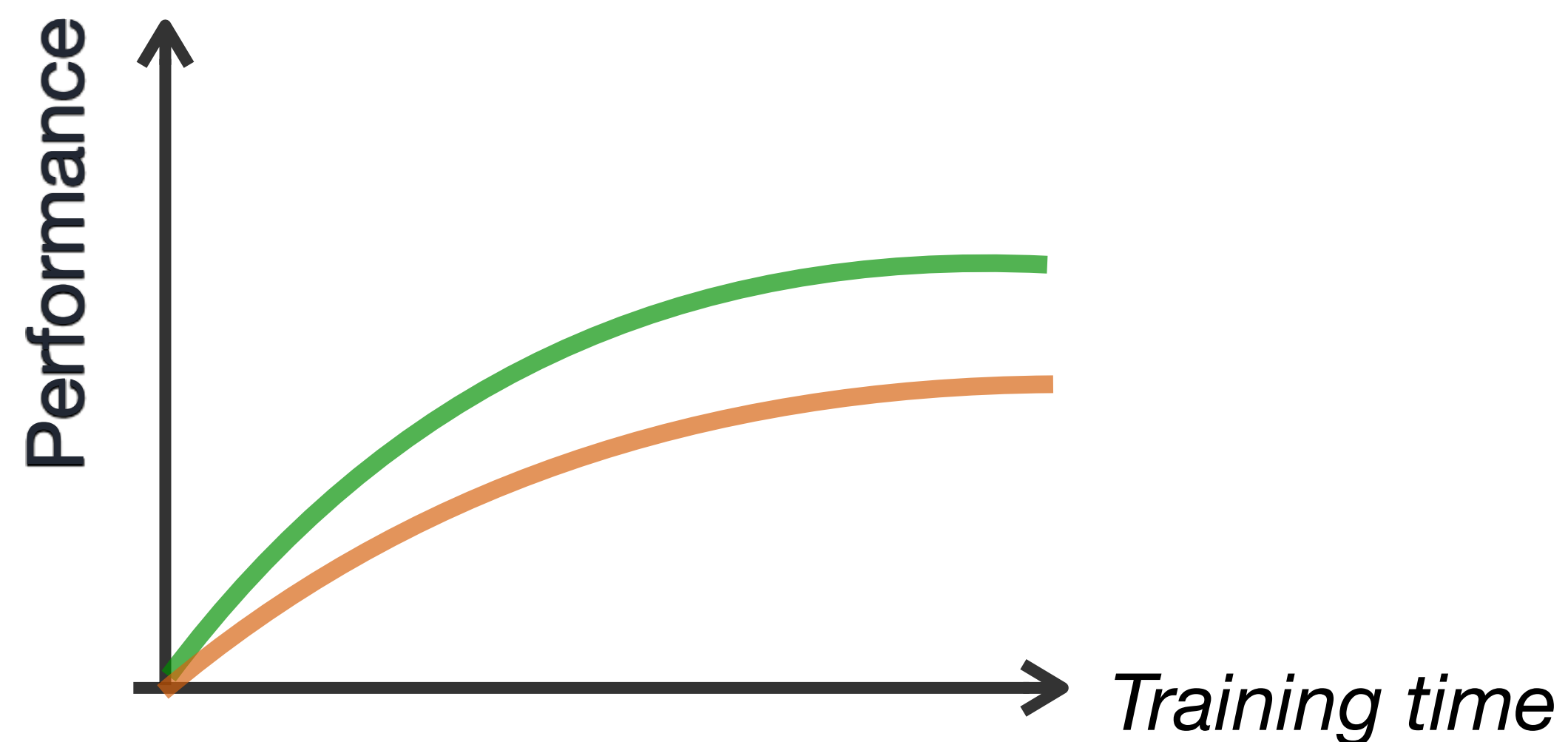
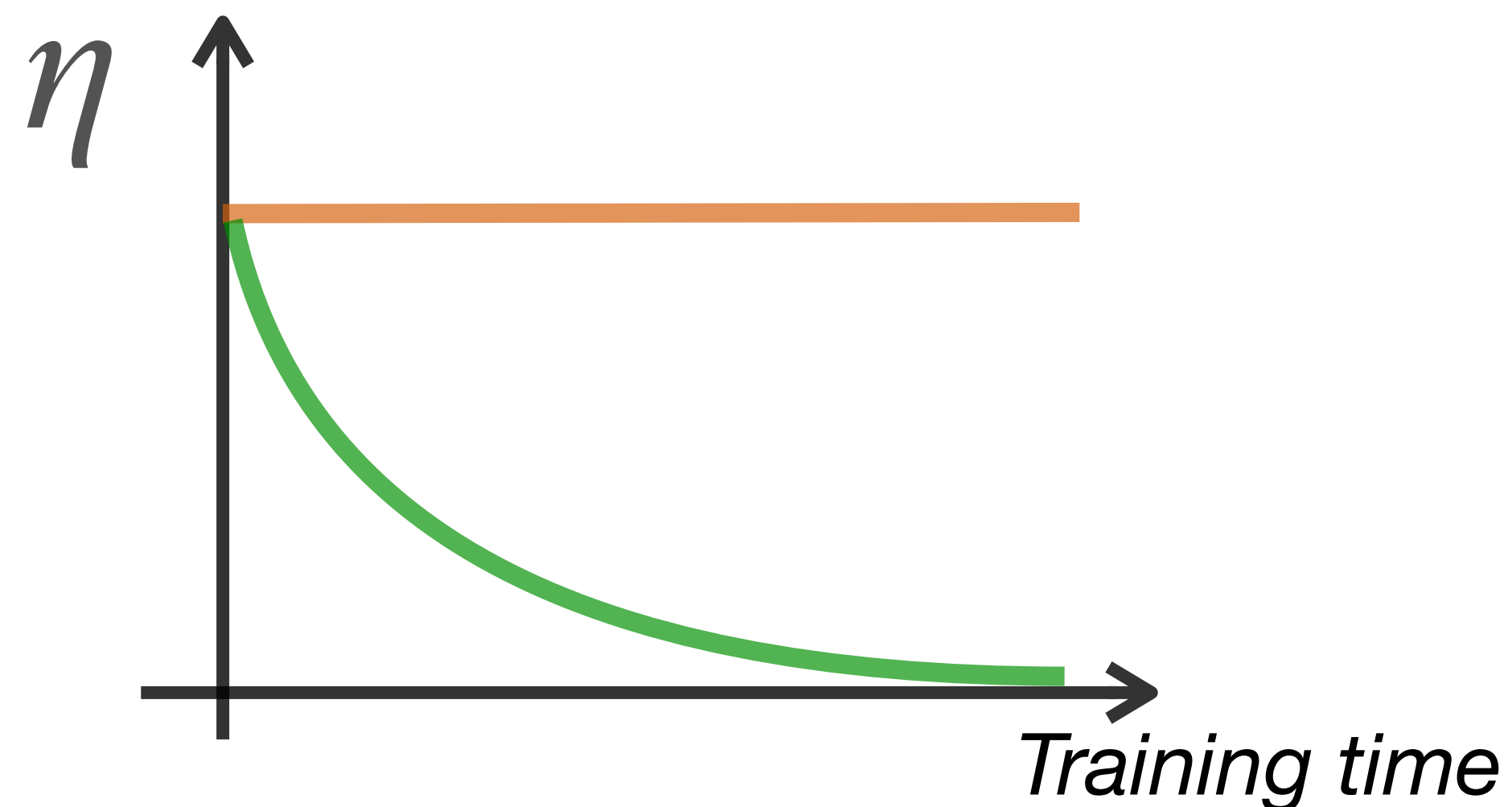


Training protocols example: learning rate

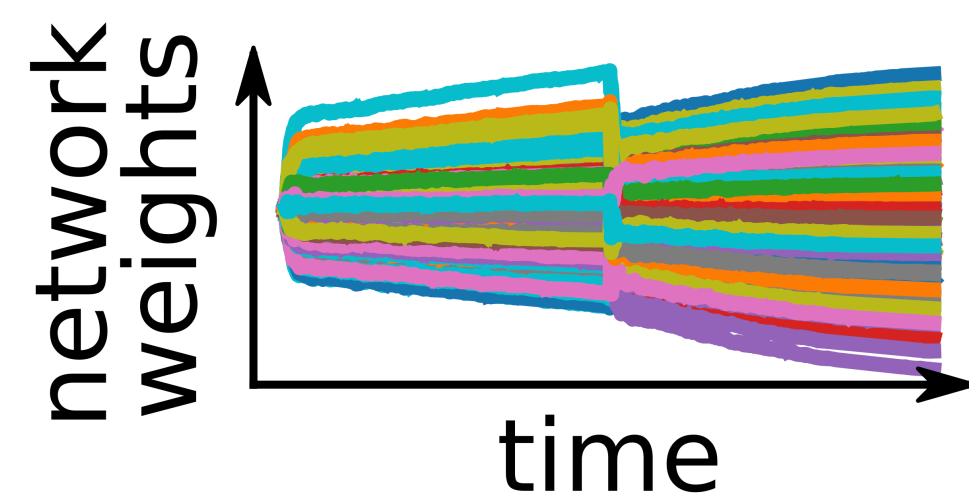


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

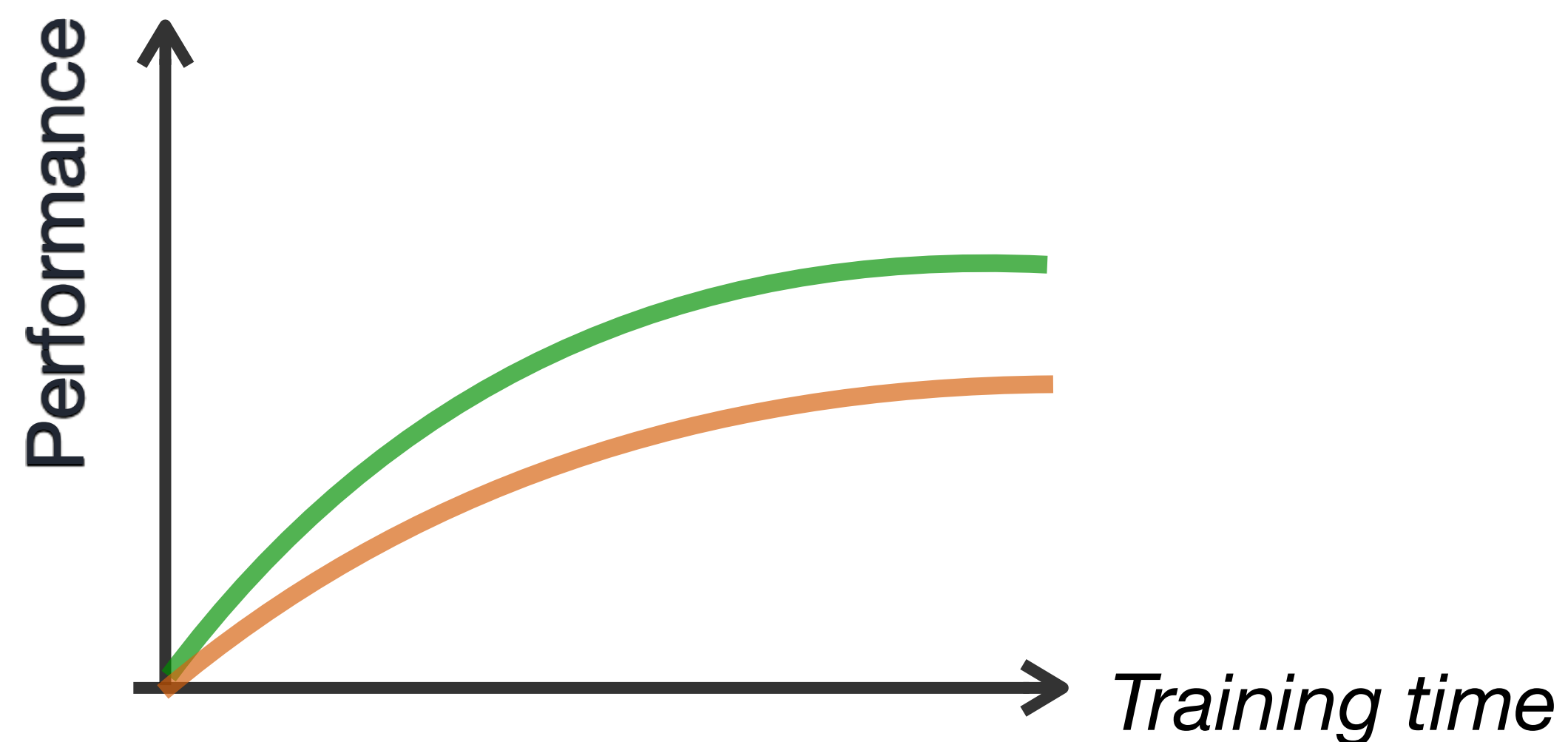
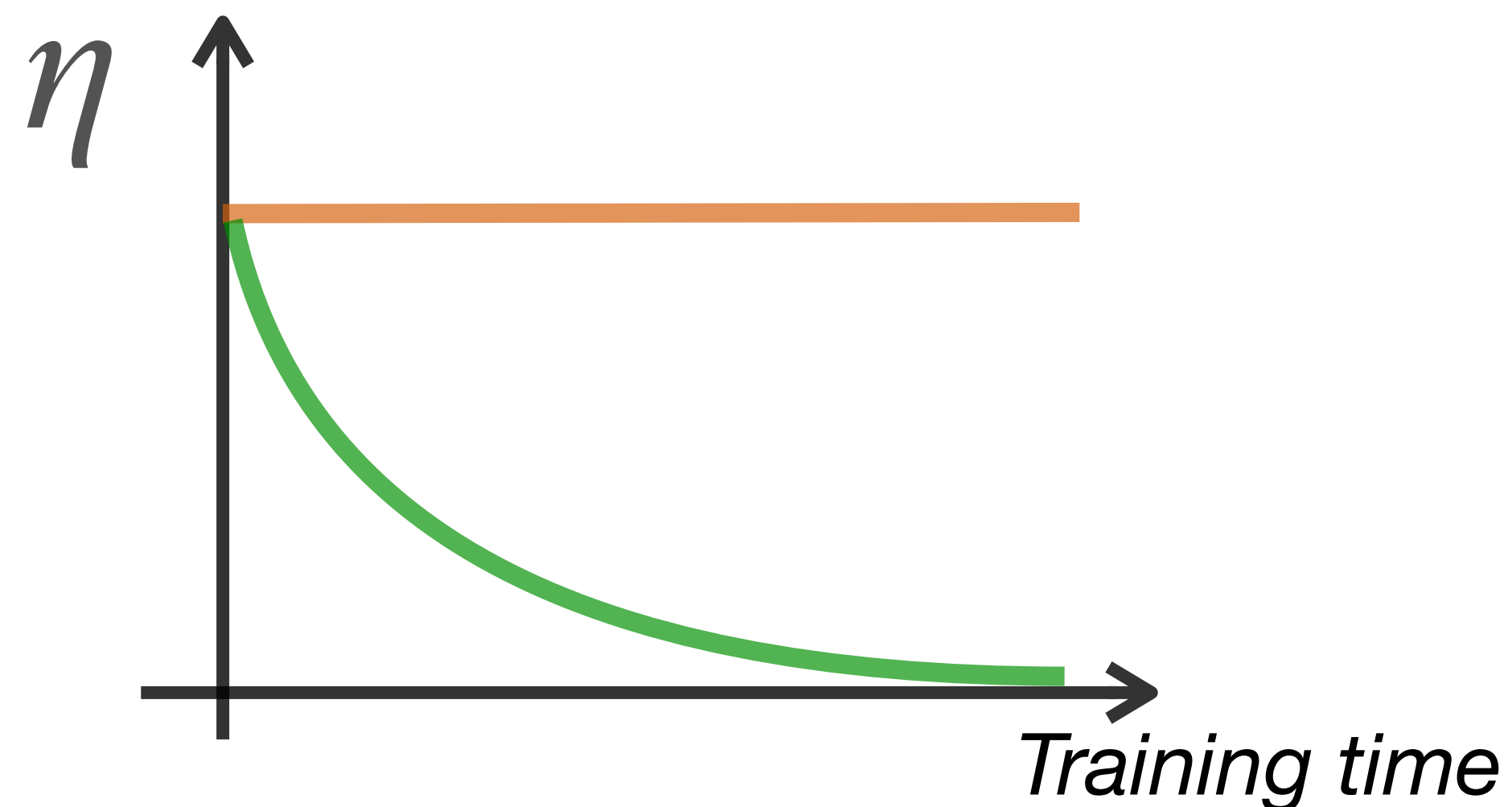


Training protocols example: learning rate



$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

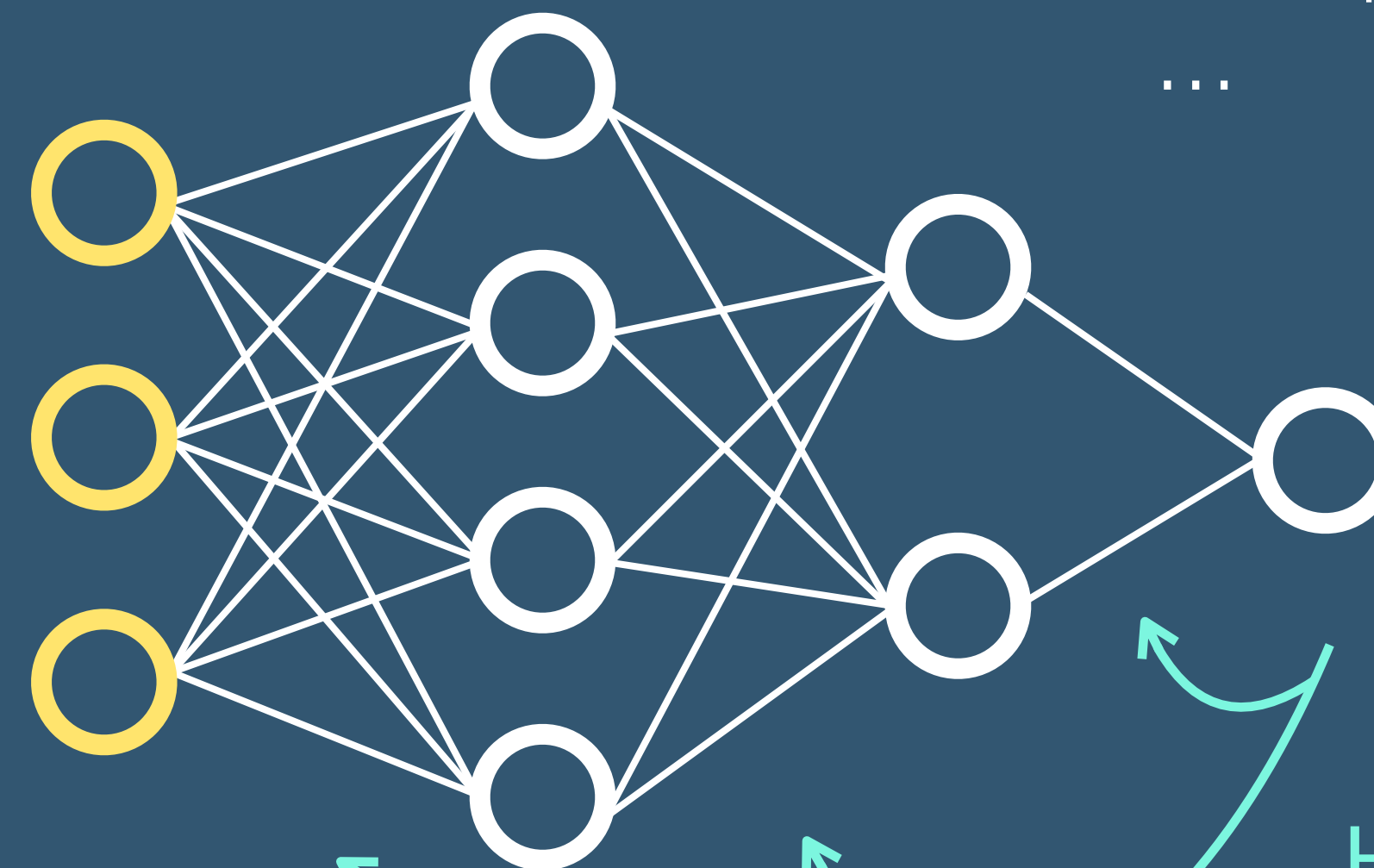


Can we compute the *optimal** strategy ?

* In terms of the final performance

Training protocols

Structured
Data / Task



Model optimization:

- Pruning
- Knowledge distillation
- Dropout

...

Dynamic data / task selection:

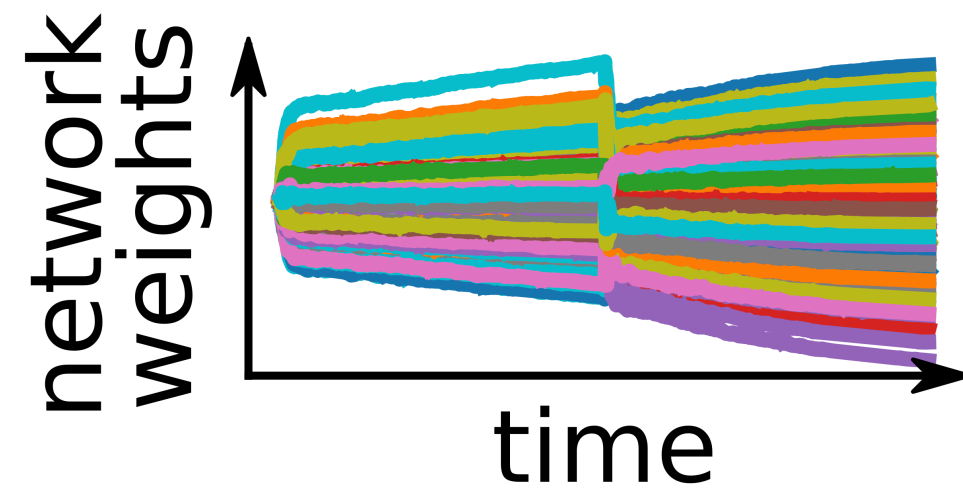
- Active learning
- Curriculum learning
- Transfer learning
- Multi-task learning
- ...

Optimization
Algorithm

Hyper-parameter schedules:

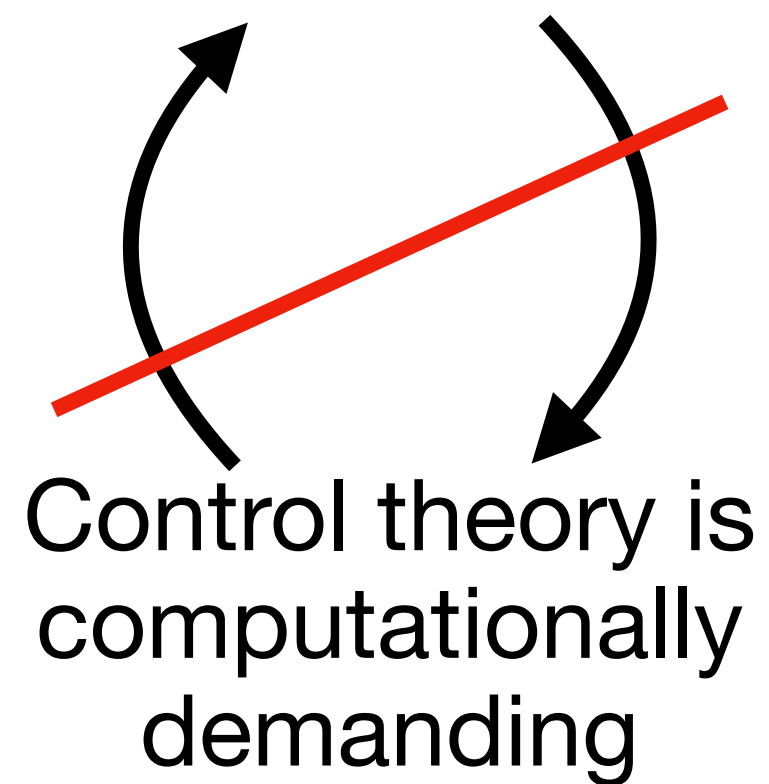
- Learning rate
- Momentum
- Batch size
- Regularization
- ...

Dimensionality reduction + optimal control

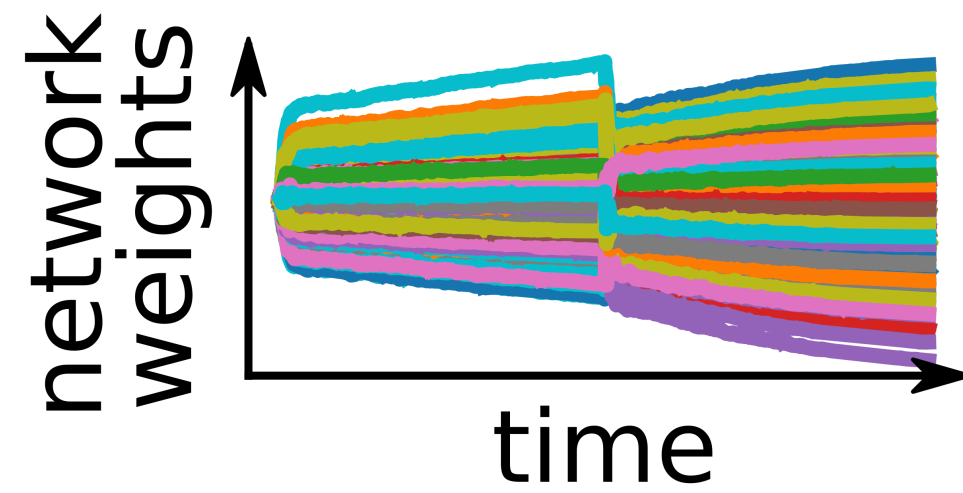


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics

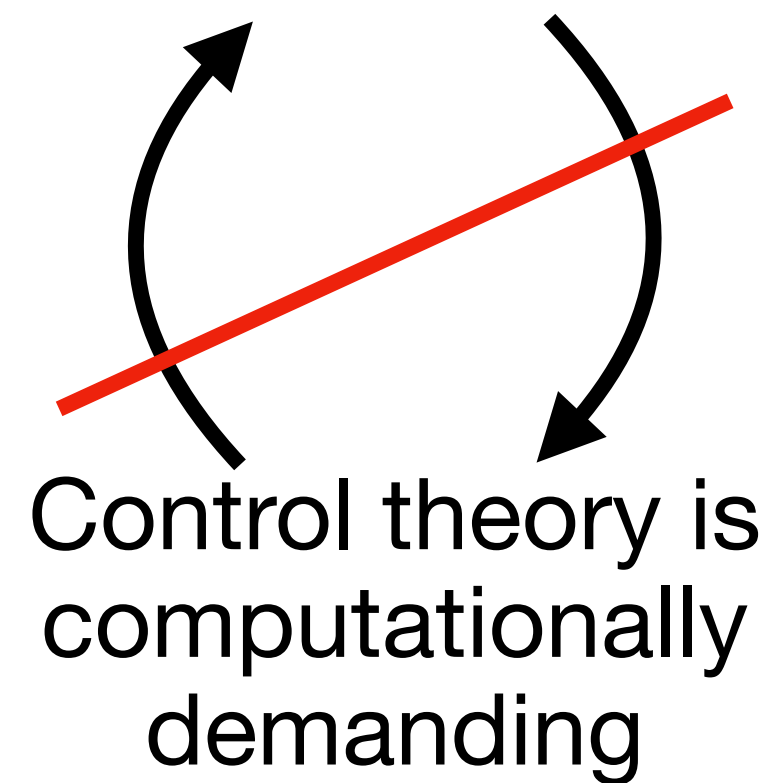


Dimensionality reduction + optimal control

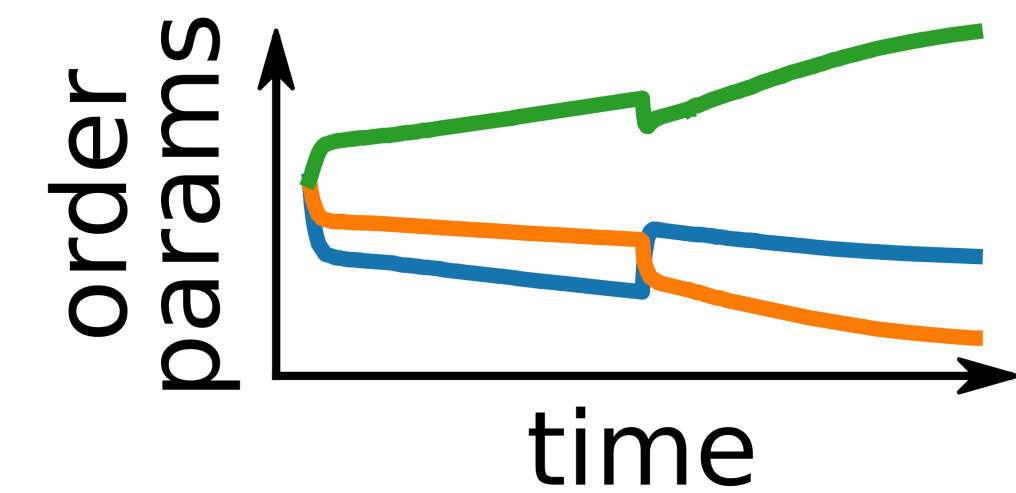


$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics



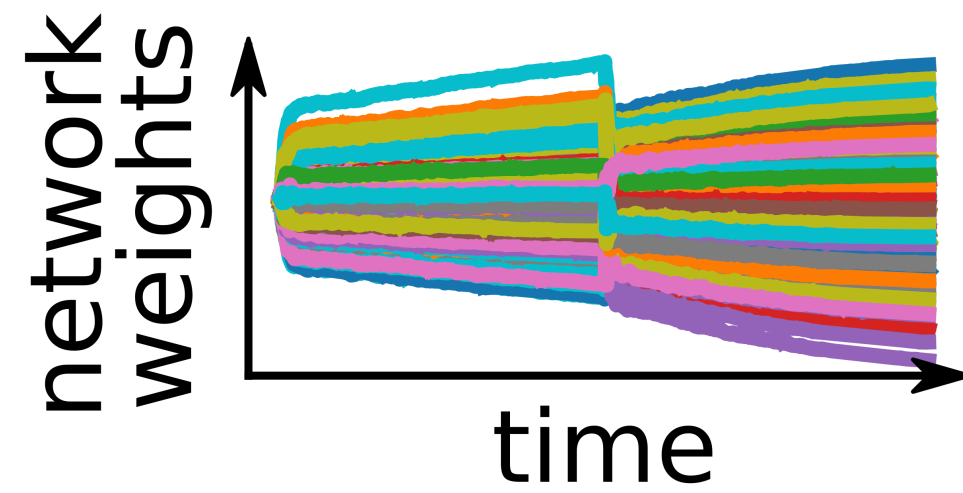
$$N \rightarrow \infty$$



$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u}) \quad \text{control variables}$$

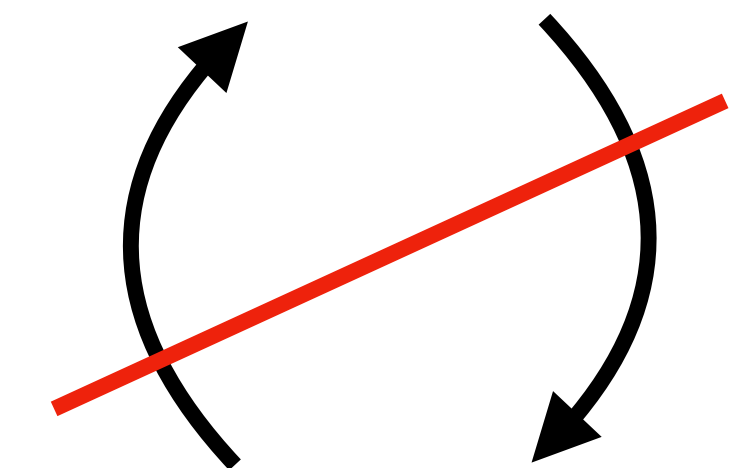
Low-dimensional
effective description

Dimensionality reduction + optimal control



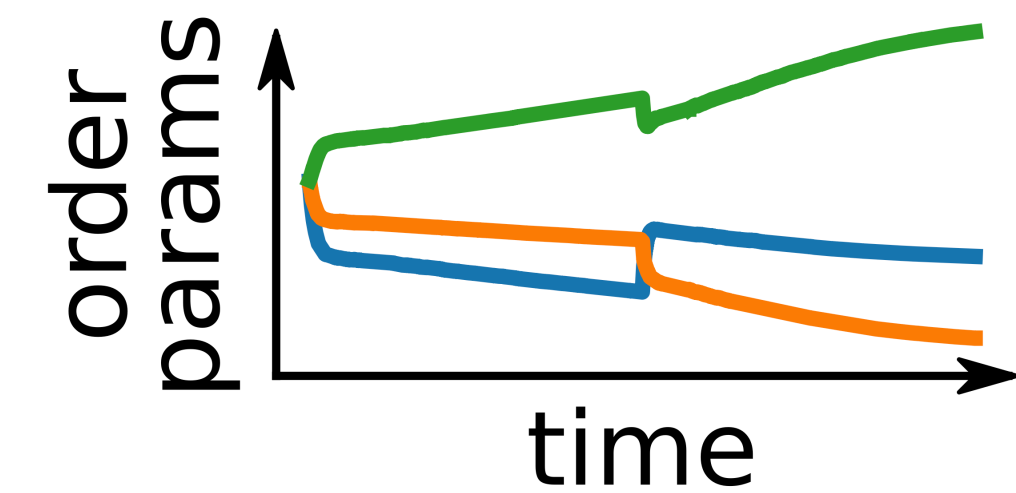
$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu}$$

High-dimensional
complex dynamics



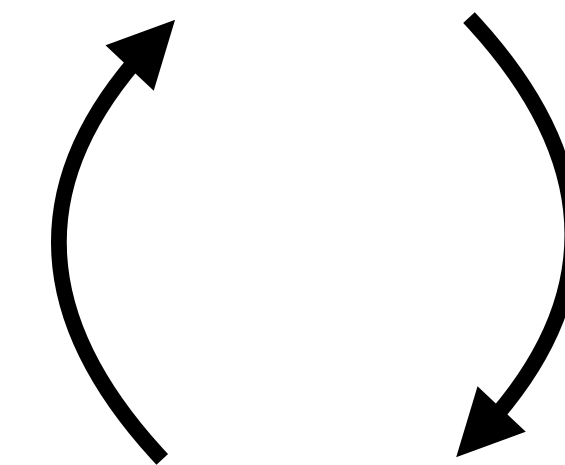
Control theory is
computationally
demanding

$$N \rightarrow \infty$$



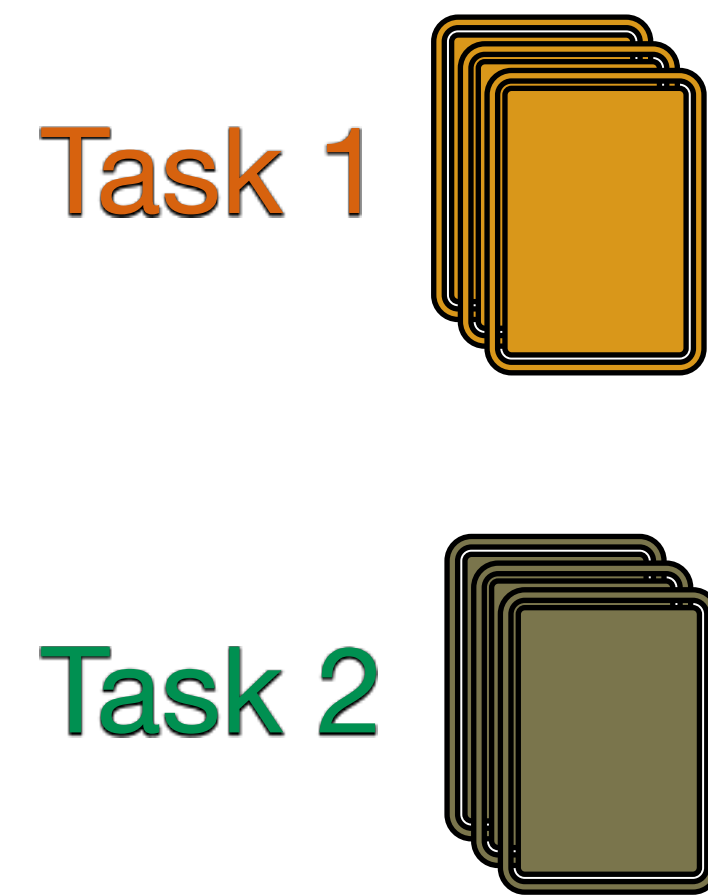
$$\dot{\mathbf{Q}} = f(\mathbf{Q}, \mathbf{u}) \quad \text{control variables}$$

Low-dimensional
effective description

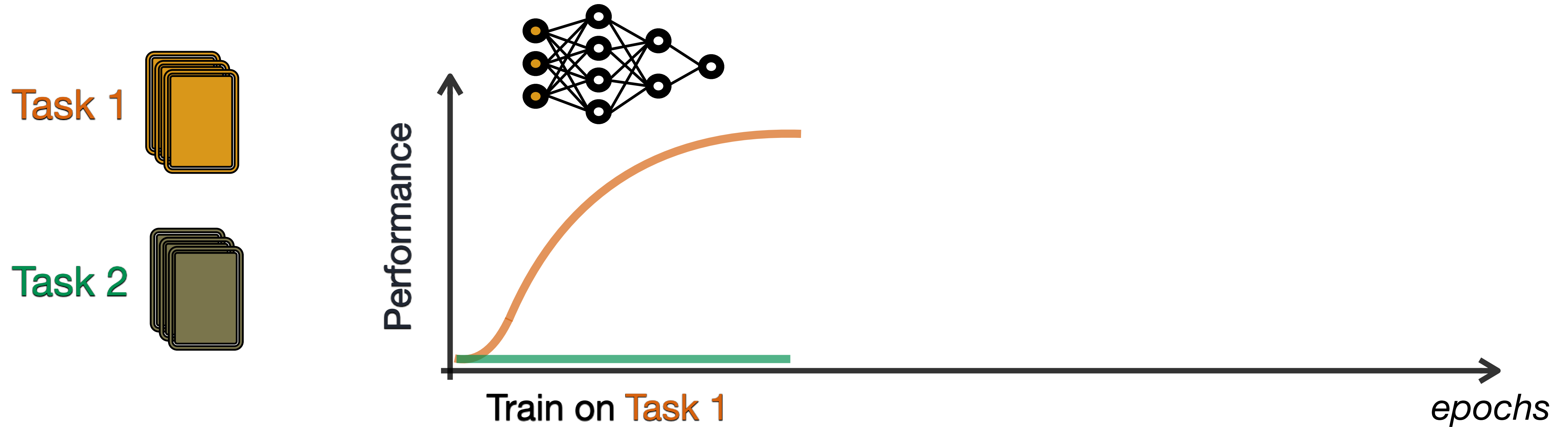


Control theory can be
applied

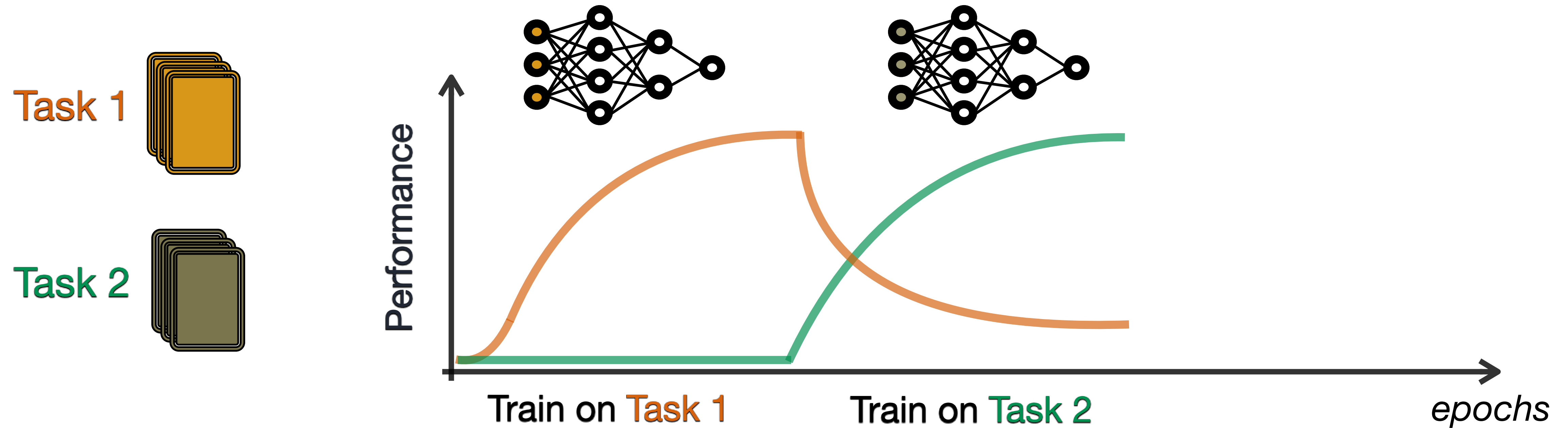
Continual learning & catastrophic forgetting



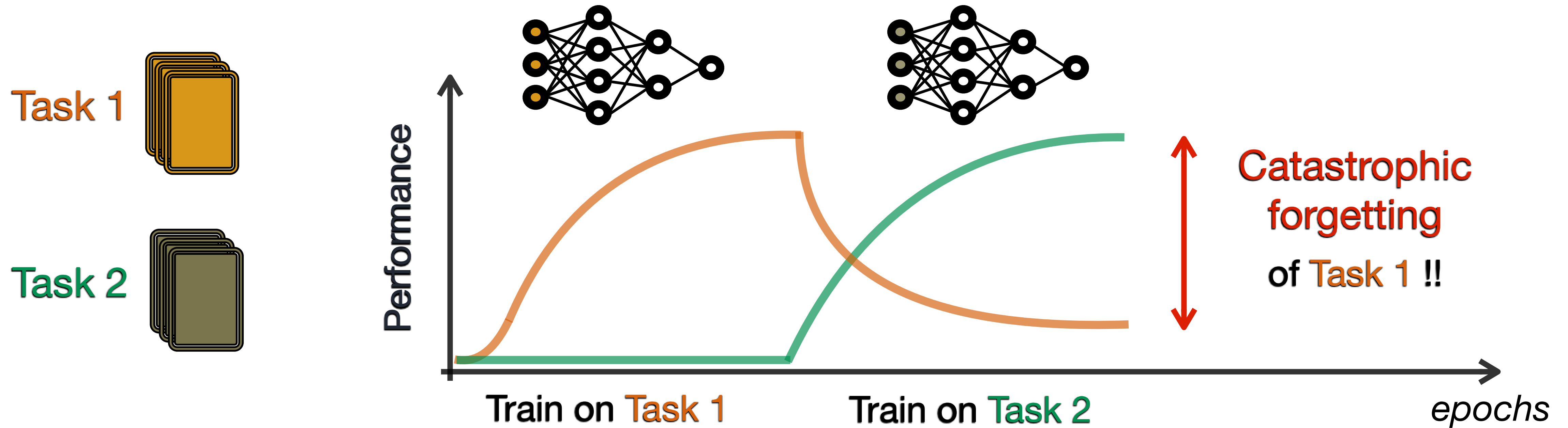
Continual learning & catastrophic forgetting



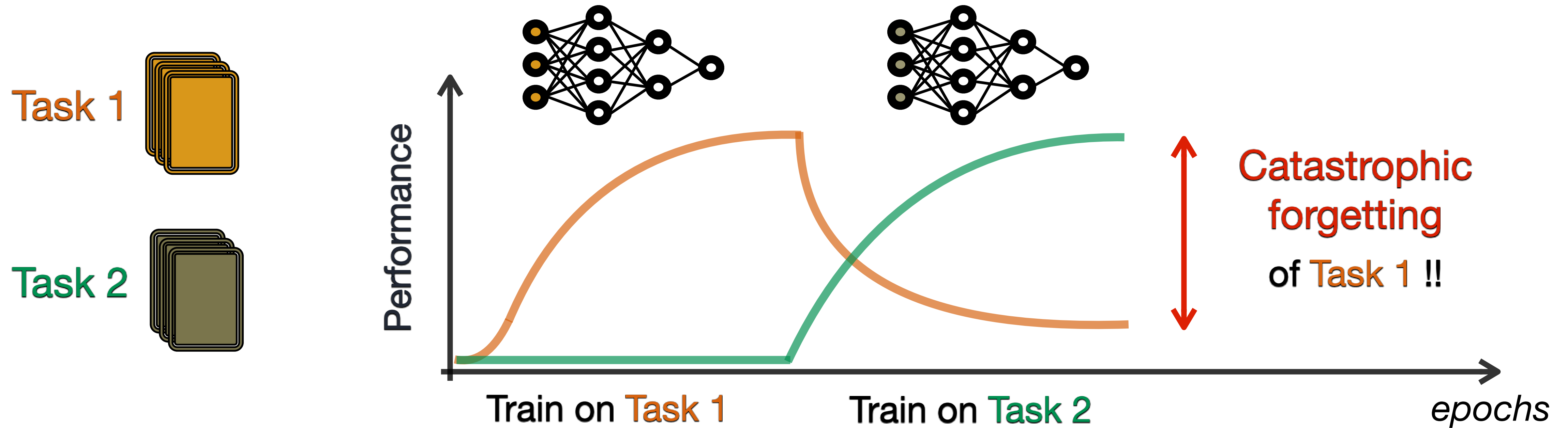
Continual learning & catastrophic forgetting



Continual learning & catastrophic forgetting



Continual learning & catastrophic forgetting



ML (empirical):

ML (theory):

A teacher-student model of continual learning

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

Task 1

$$\mathcal{D}_1 = \{x_i^{(1)}, y_i^{(1)}\}$$

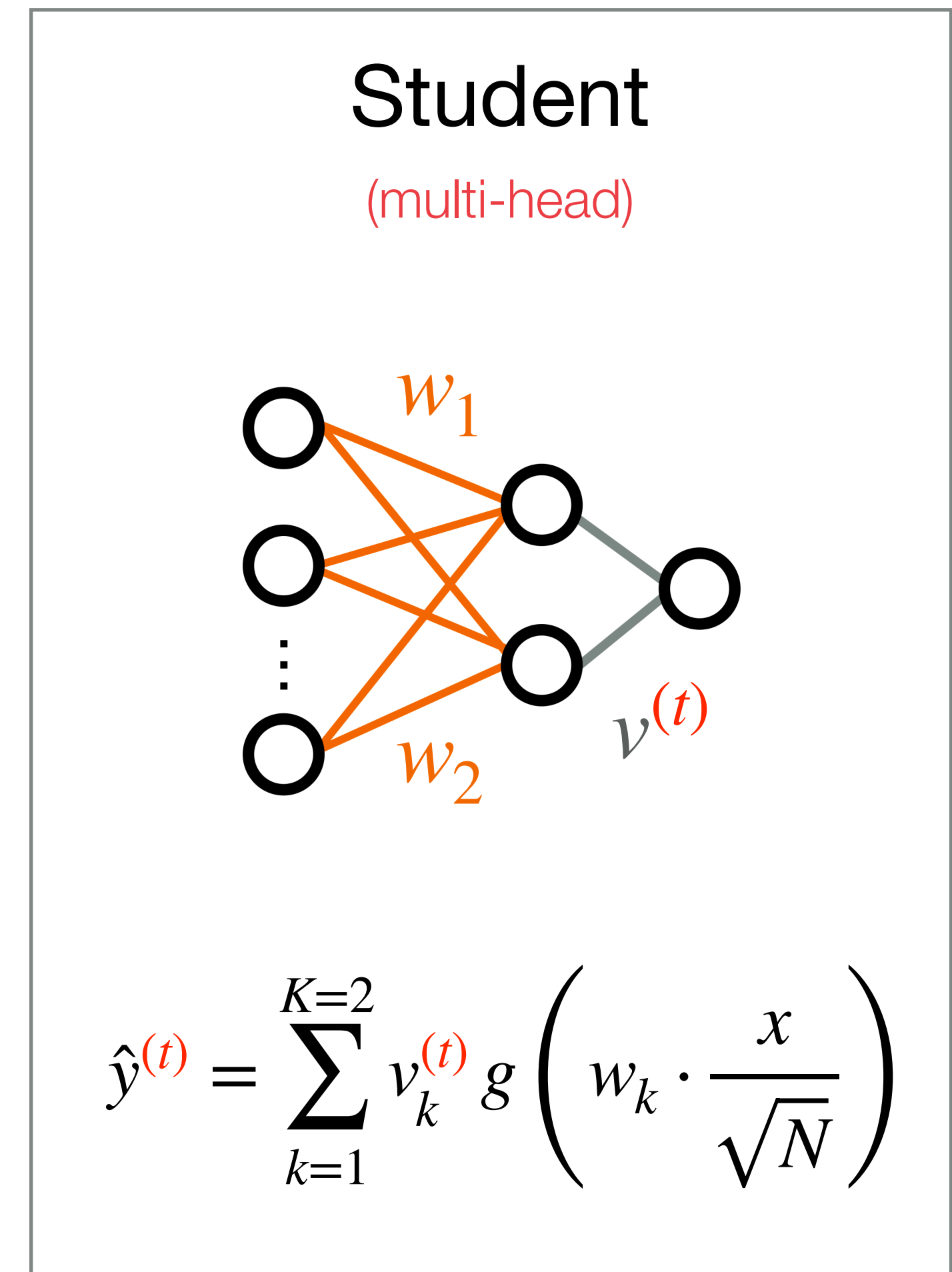
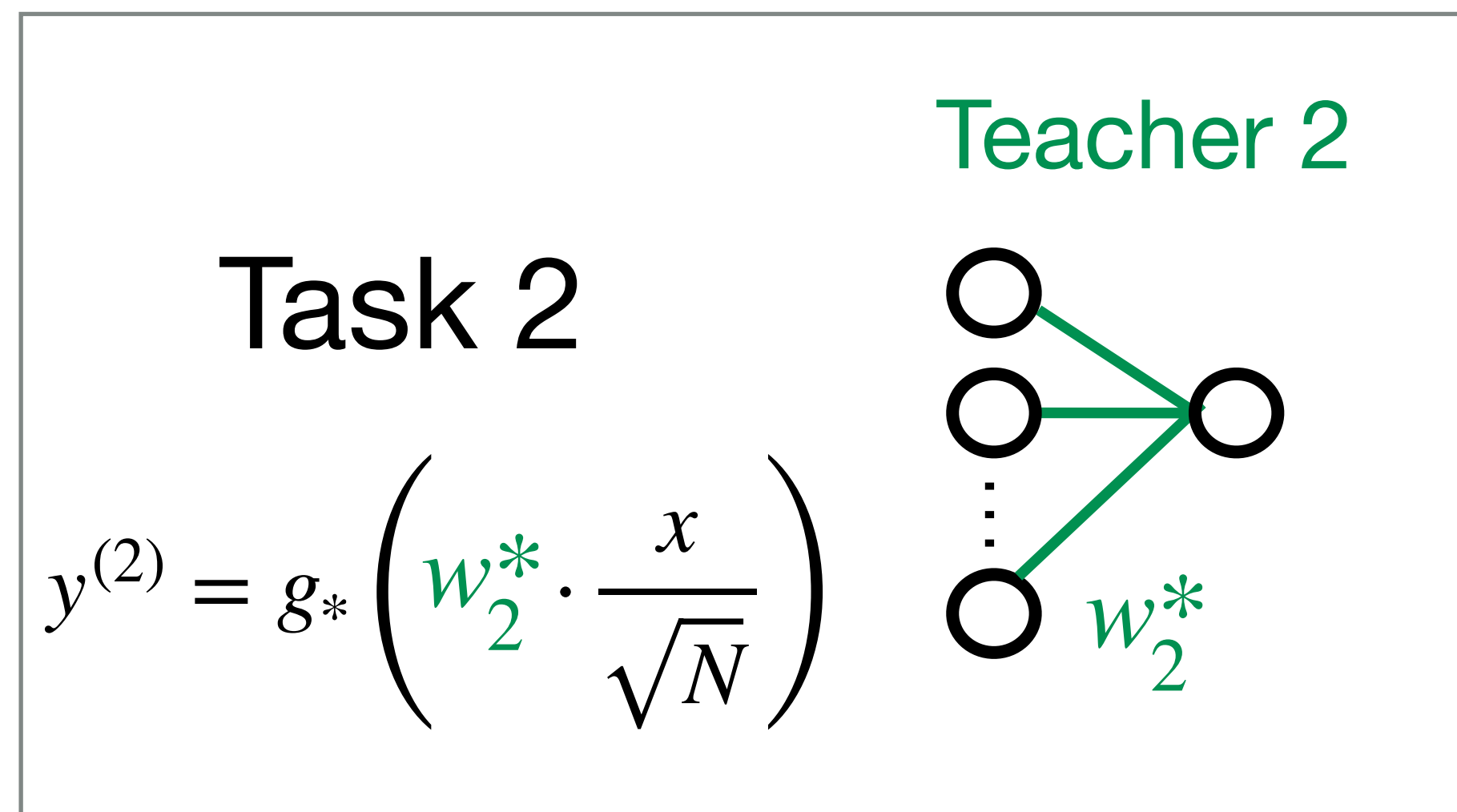
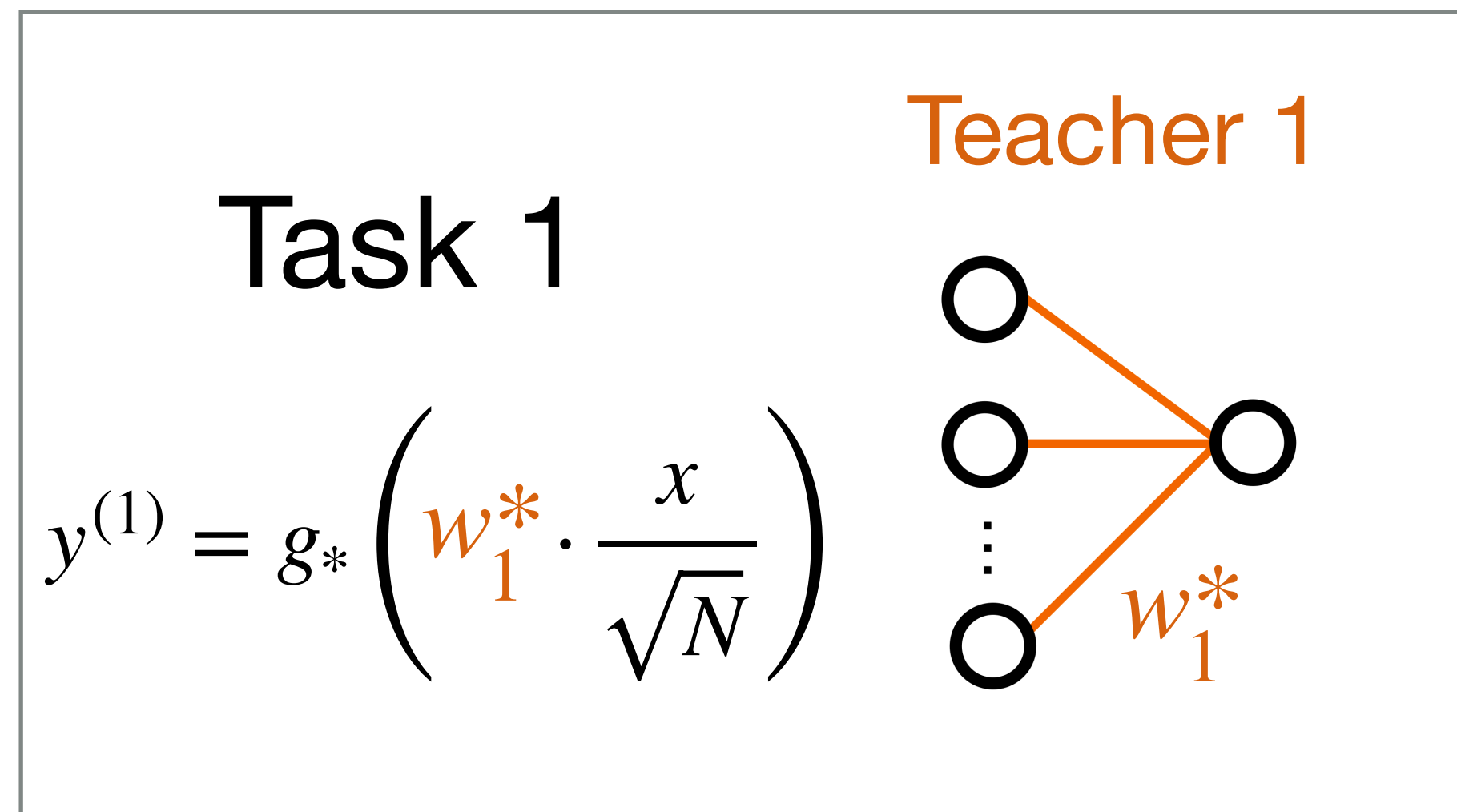
$$x \sim \mathcal{N}(0,1) \in \mathbf{R}^N \quad N \gg 1$$

Task 2

$$\mathcal{D}_2 = \{x_i^{(2)}, y_i^{(2)}\}$$

A teacher-student model of continual learning

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



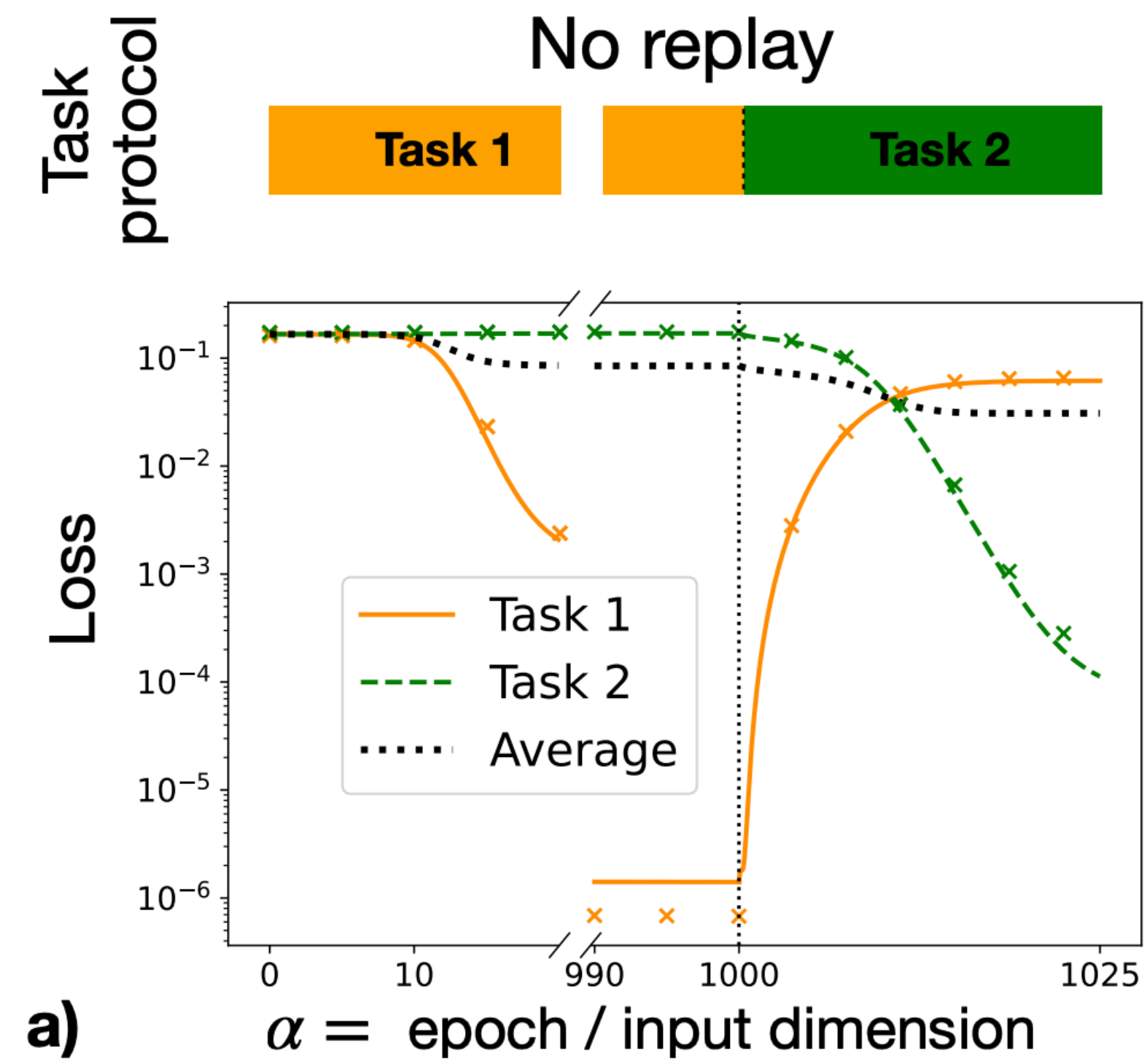
A teacher-student model of continual learning

$$\mathbf{w}^{\mu+1} = \mathbf{w}^{\mu} - \eta \nabla_{\mathbf{w}} \mathcal{L}^{\mu} \xrightarrow{N \rightarrow \infty} \frac{dQ(\alpha)}{d\alpha} = f_Q(Q(\alpha), \mathbf{u}(\alpha))$$

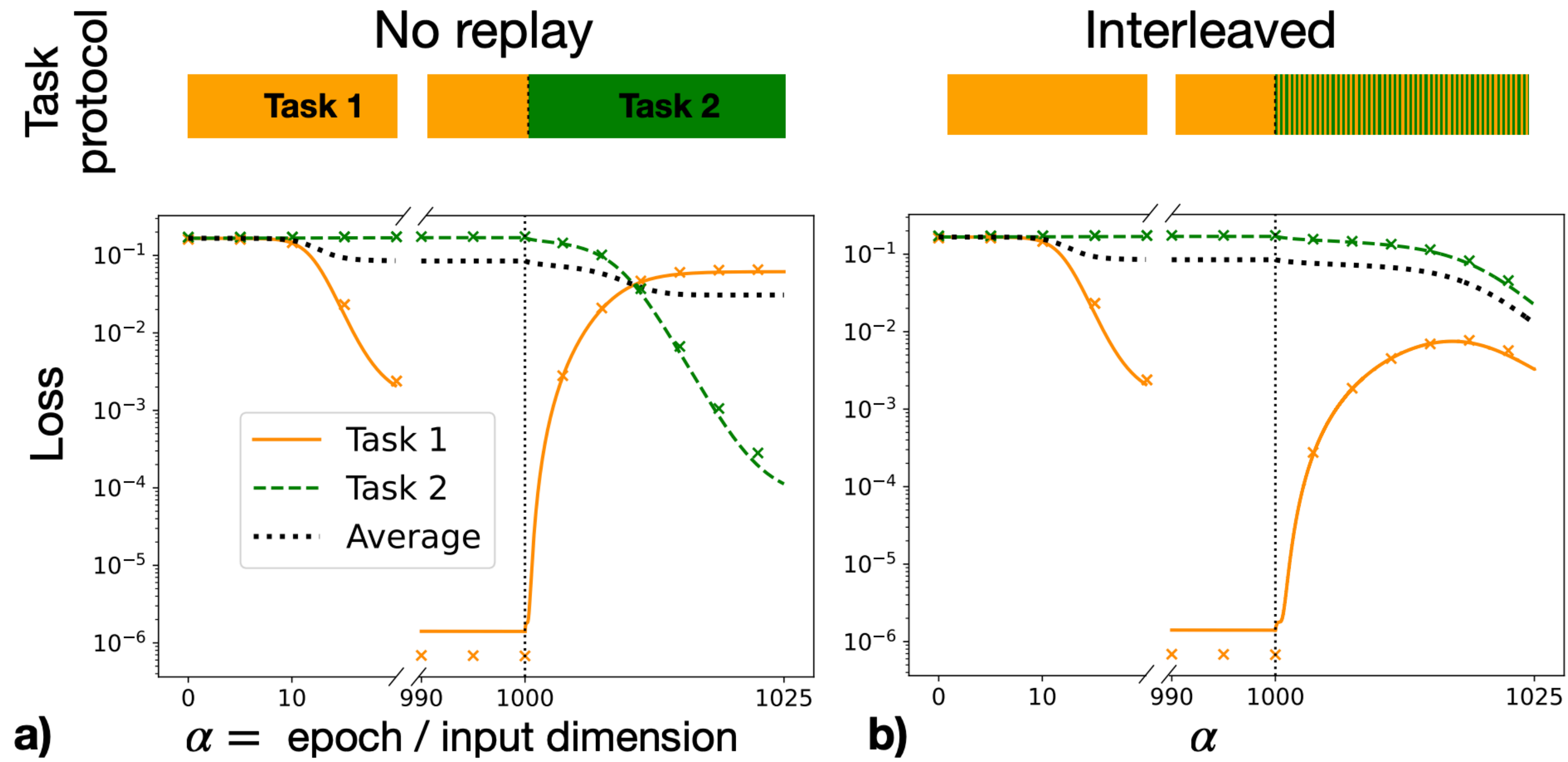
$\alpha = \mu/N \quad \alpha \in (0, \alpha_F]$

Example: "Magnetisation" $M_{kt} = \frac{w_k \cdot w_k^t}{N}$

Results: optimal strategy vs benchmarks

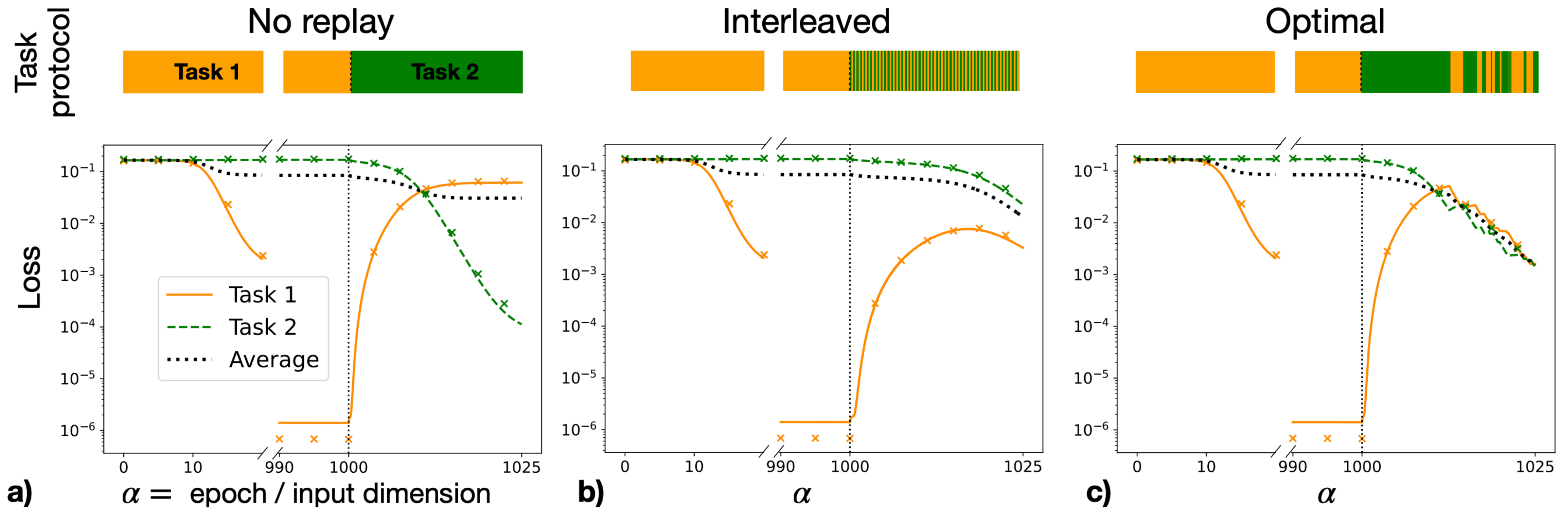


Results: optimal strategy vs benchmarks



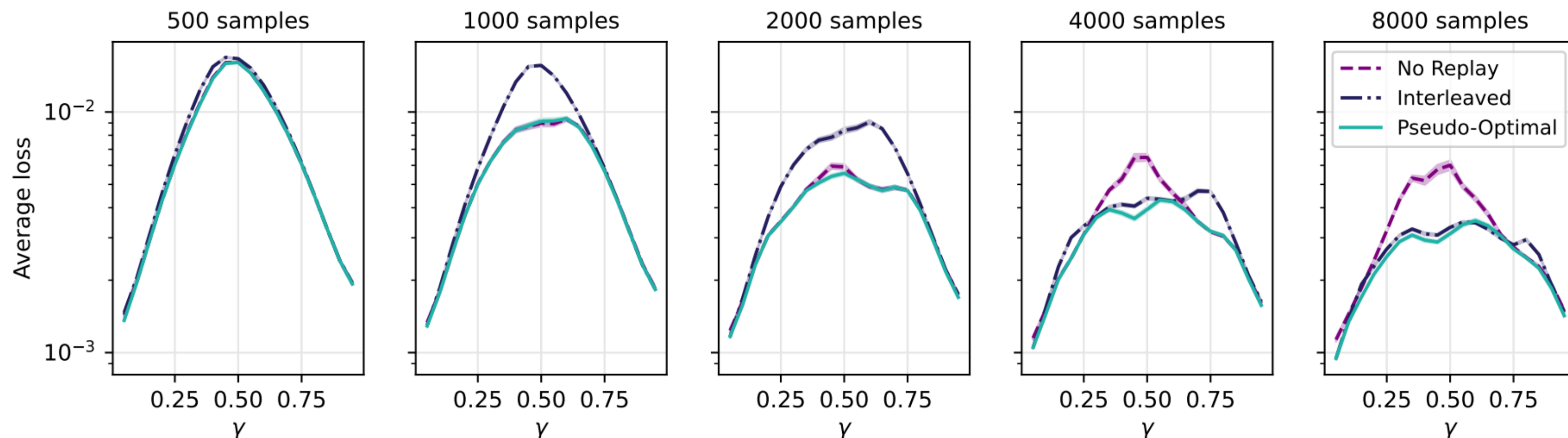
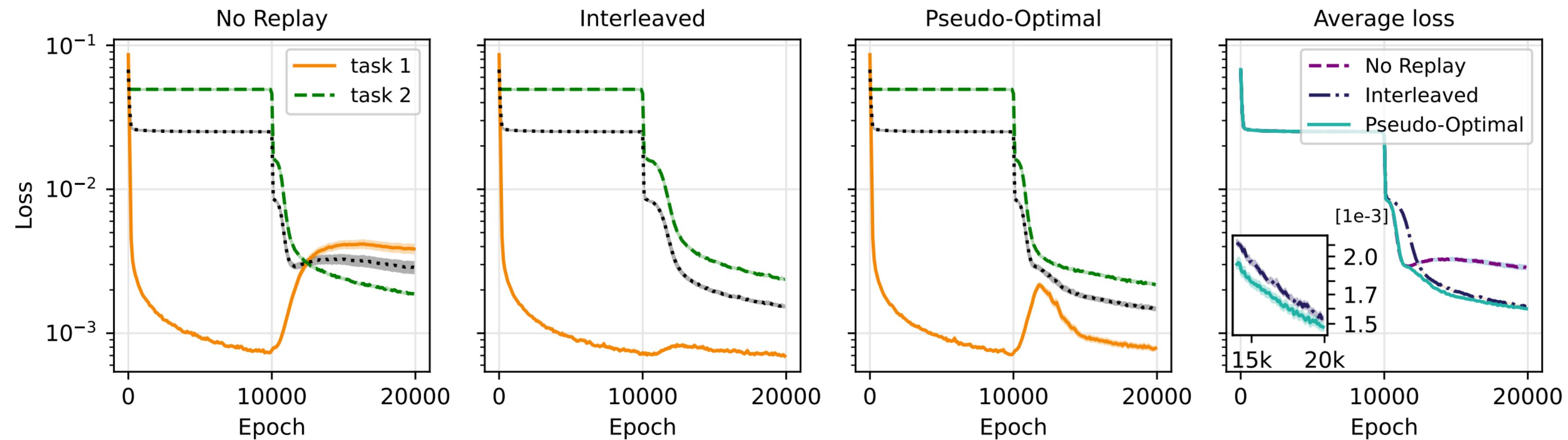
Results: optimal strategy vs benchmarks

control = task

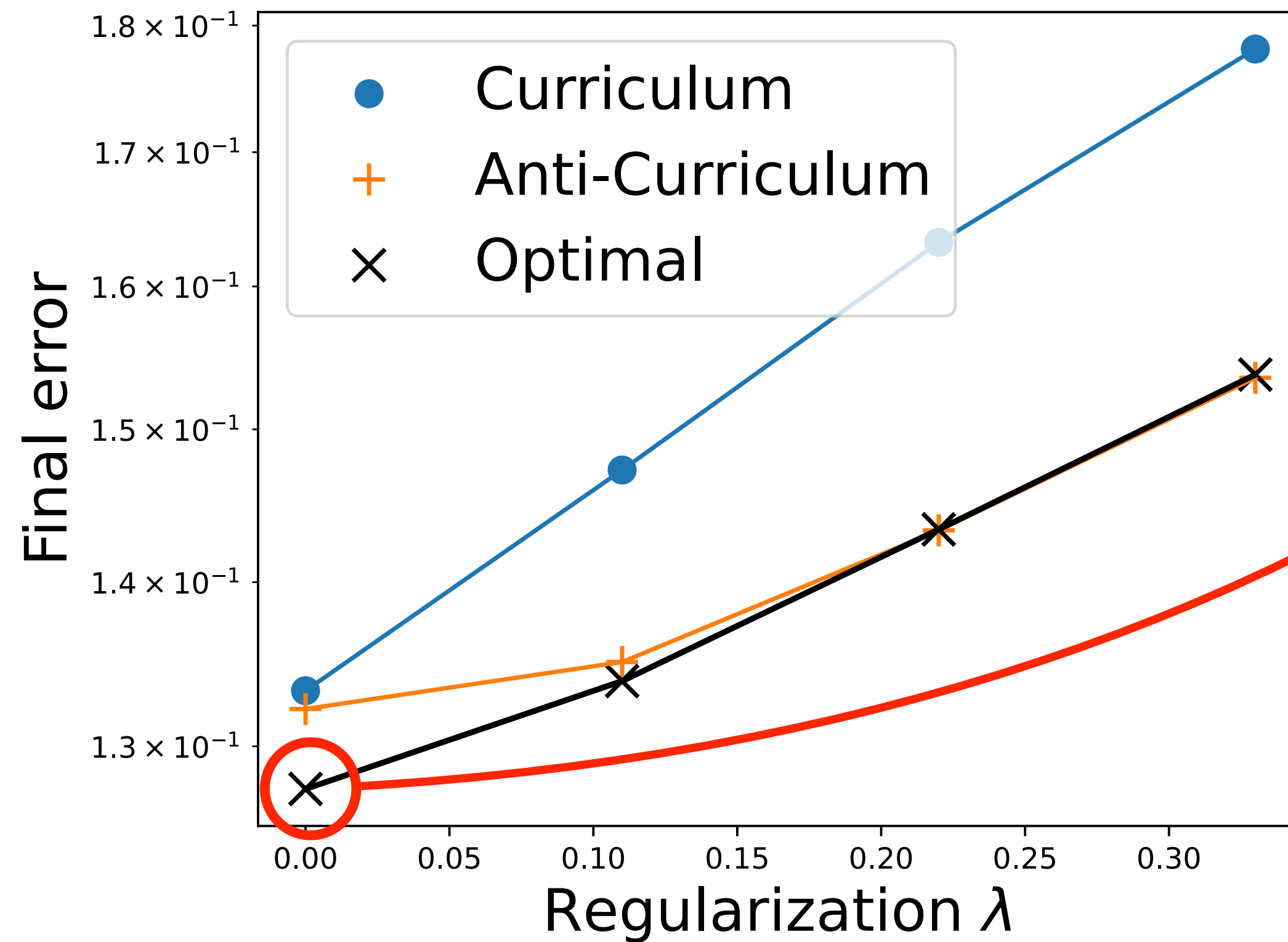


Experiments on Fashion MNIST

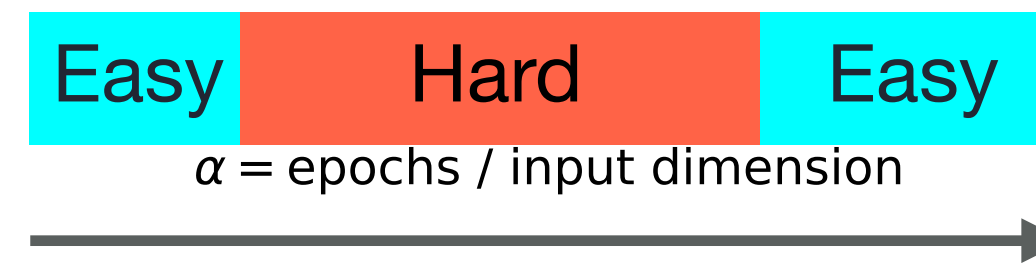
$$\mathcal{D}_1 = \{\mathbf{x}_i^{(1)}, y_i^{(1)}\}_i \quad \mathcal{D}_2 = \{\mathbf{x}_i^{(2)}, y_i^{(2)}\}_i = \{\gamma \mathbf{x}_i^{(1)} + (1 - \gamma) \tilde{\mathbf{x}}_i, \gamma y_i^{(1)} + (1 - \gamma) \tilde{y}_i\}_i$$



Optimal curriculum learning

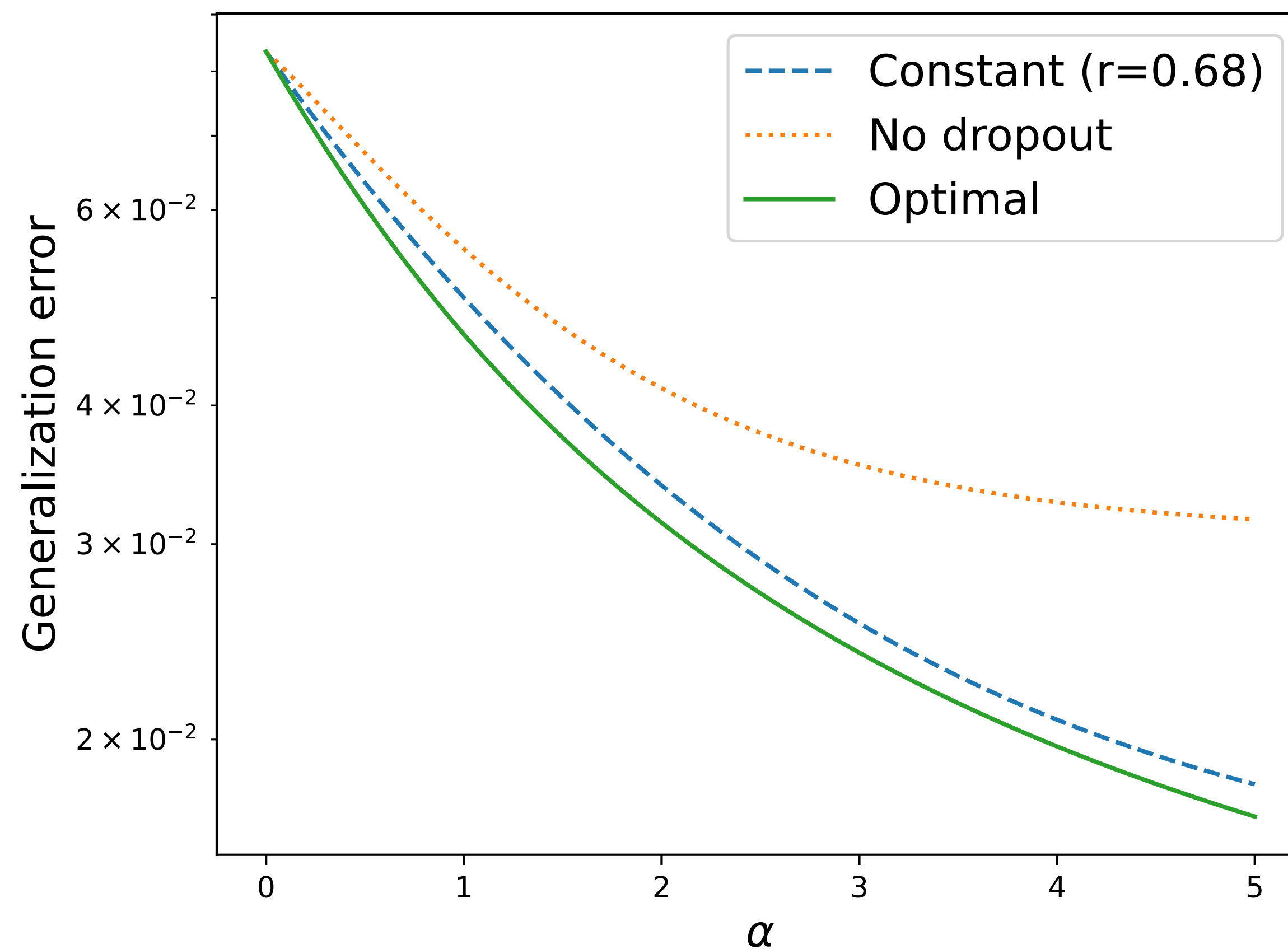
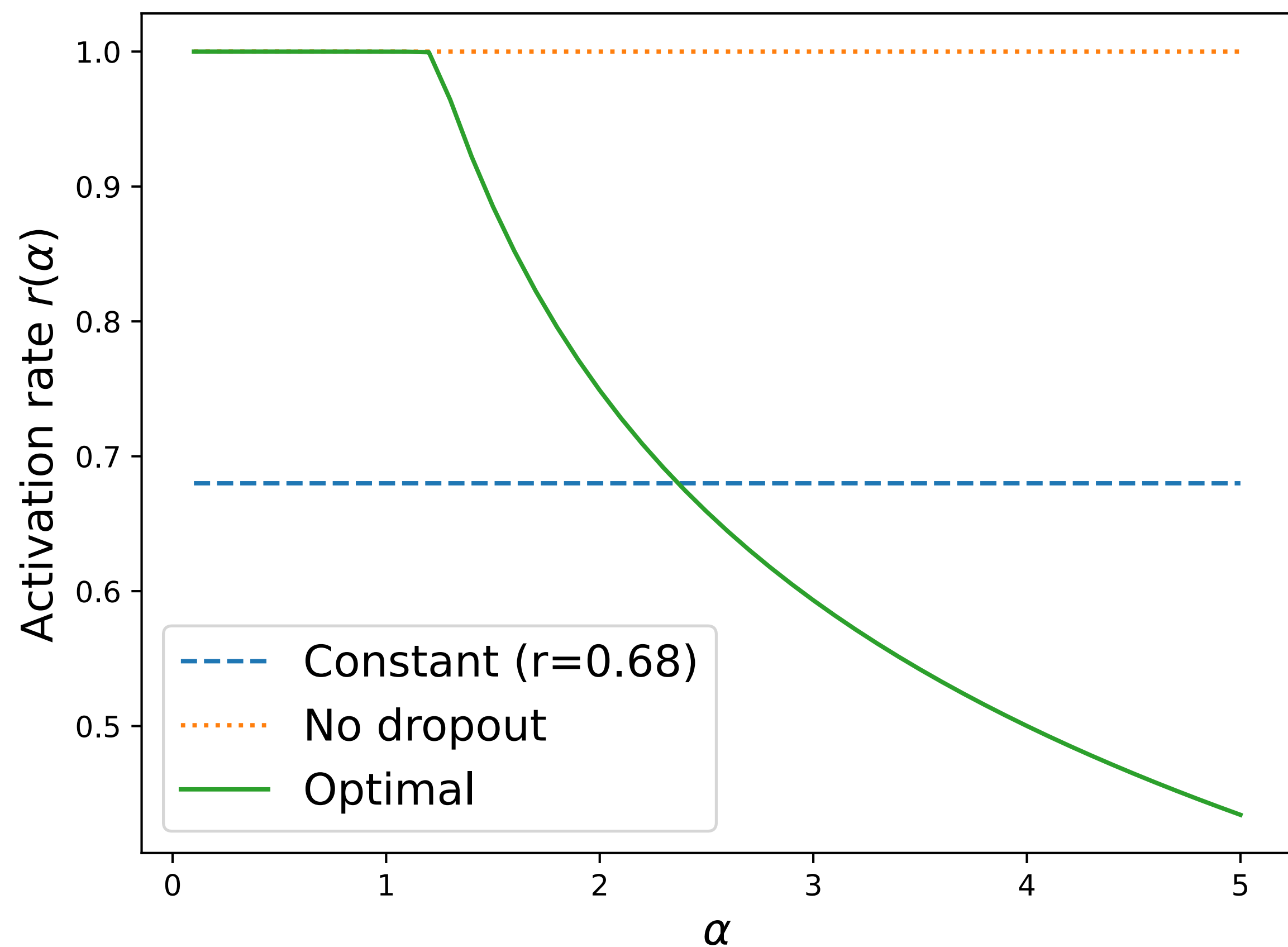


Non-monotonic
curriculum is optimal:



Conclusions & Perspectives

Optimal dropout



Conclusions & Perspectives

In summary:

Many open directions!

Thank you!

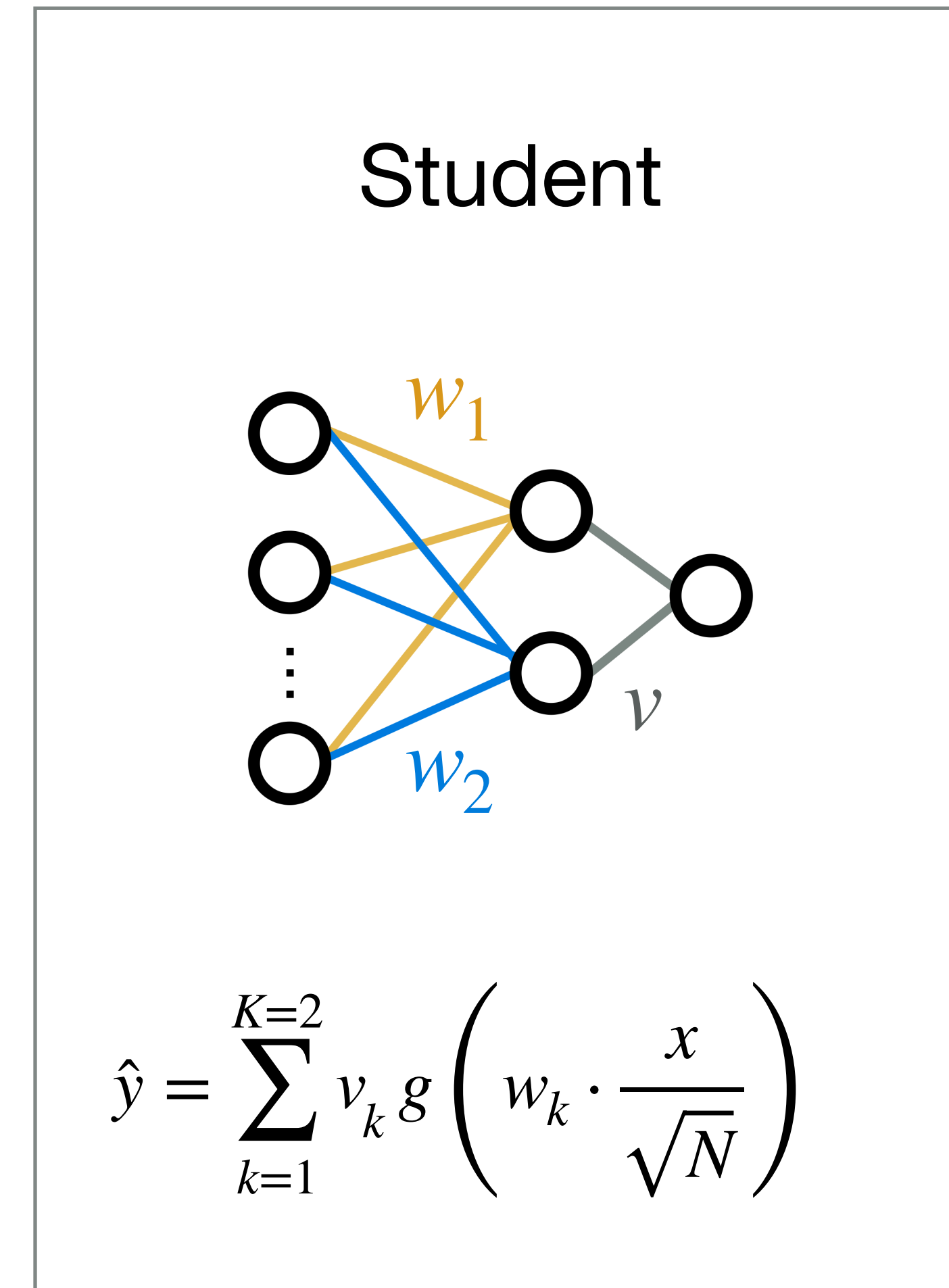
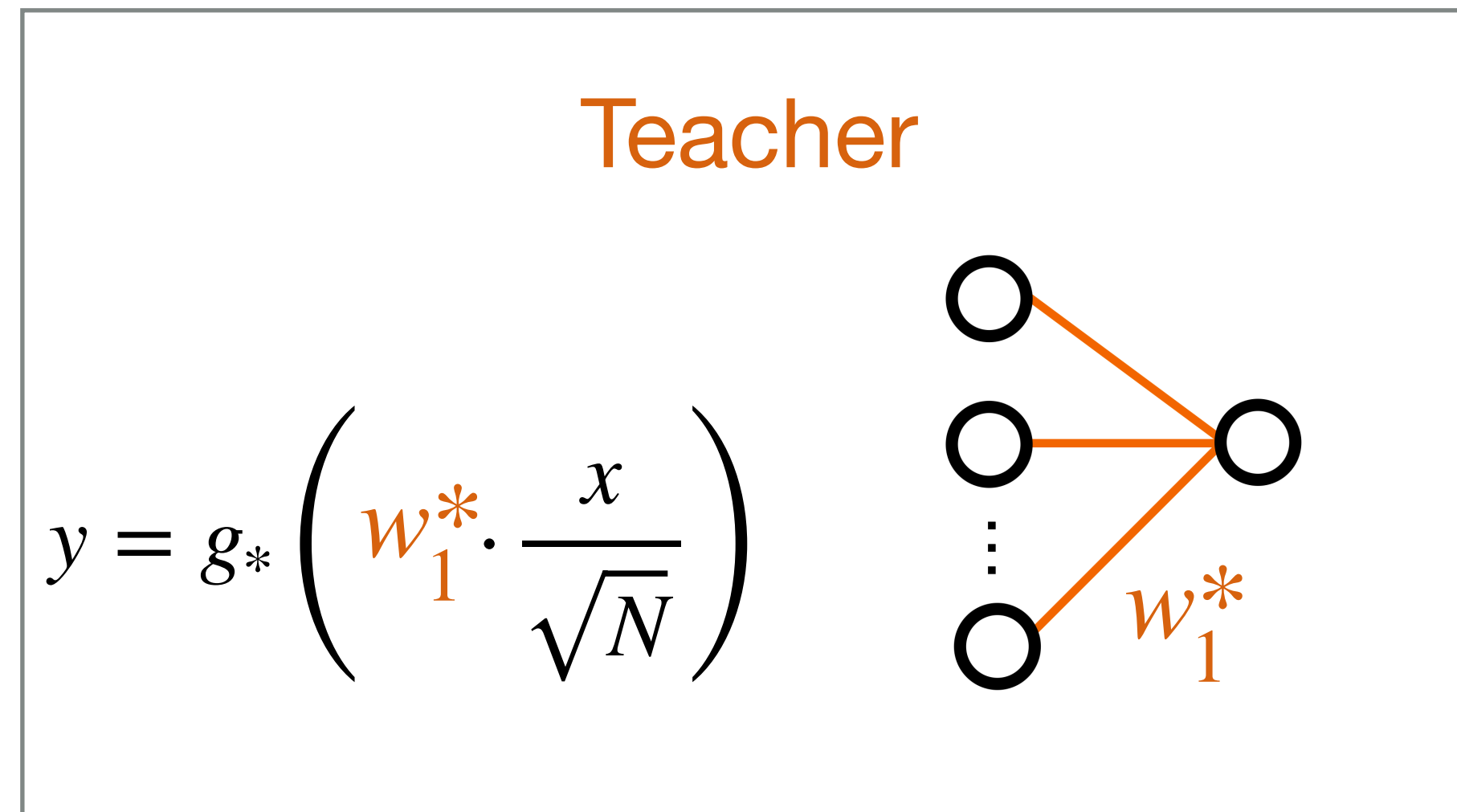
Bonus Slides

Dropout

(in progress)

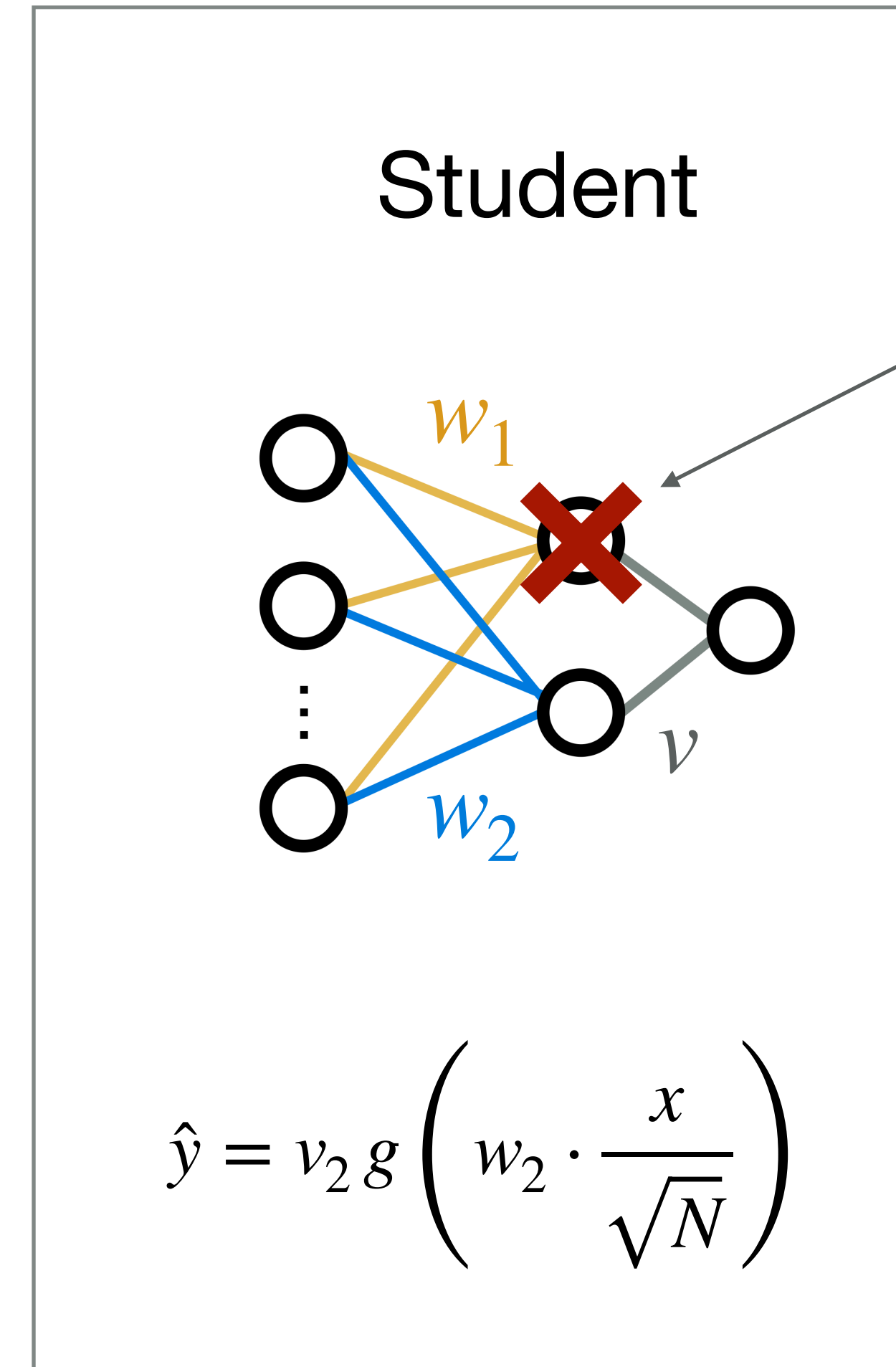
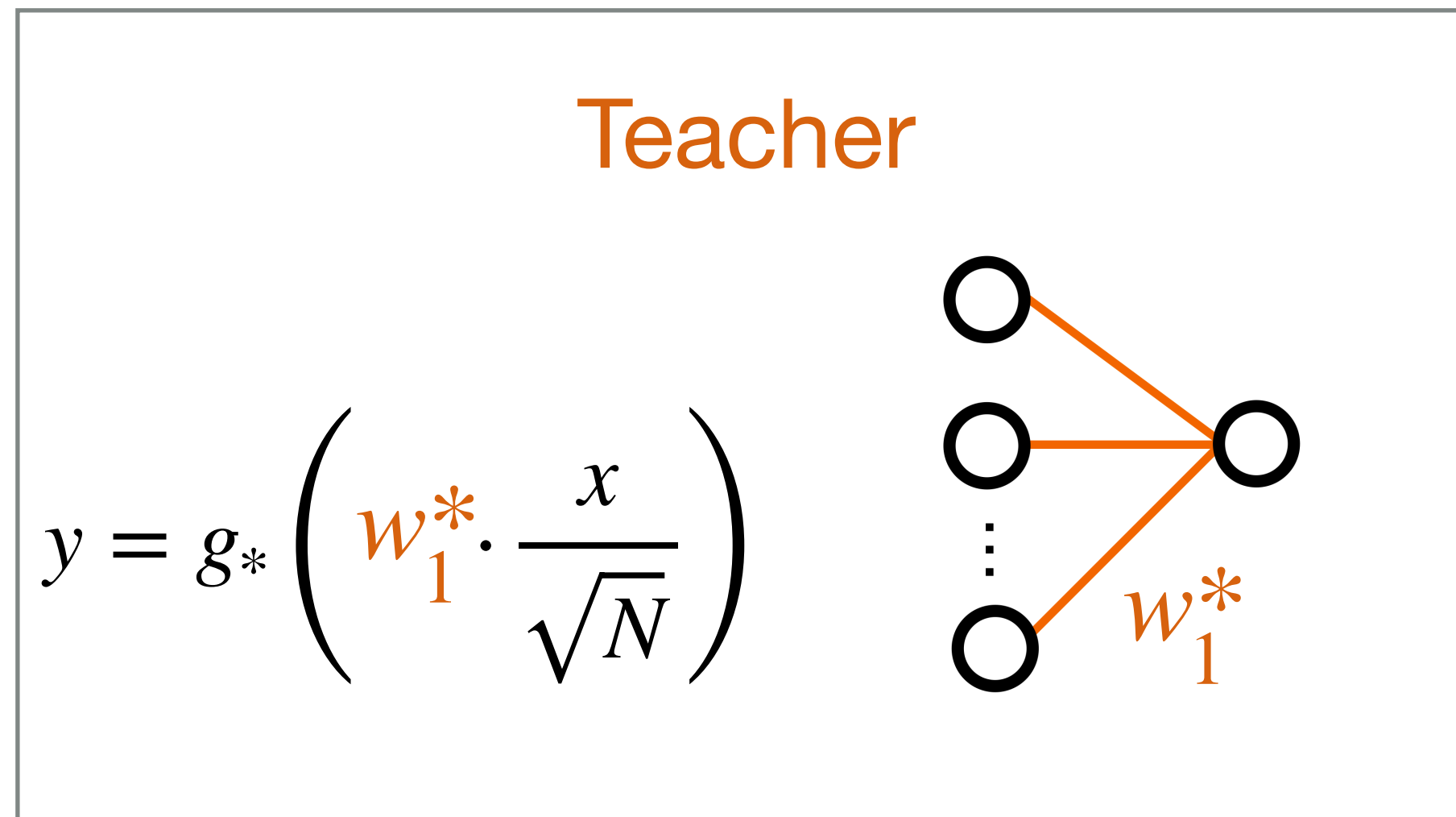
A teacher-student model of dropout

N. Srivastava, G. Hinton, et al., J. ML Res. (2014)



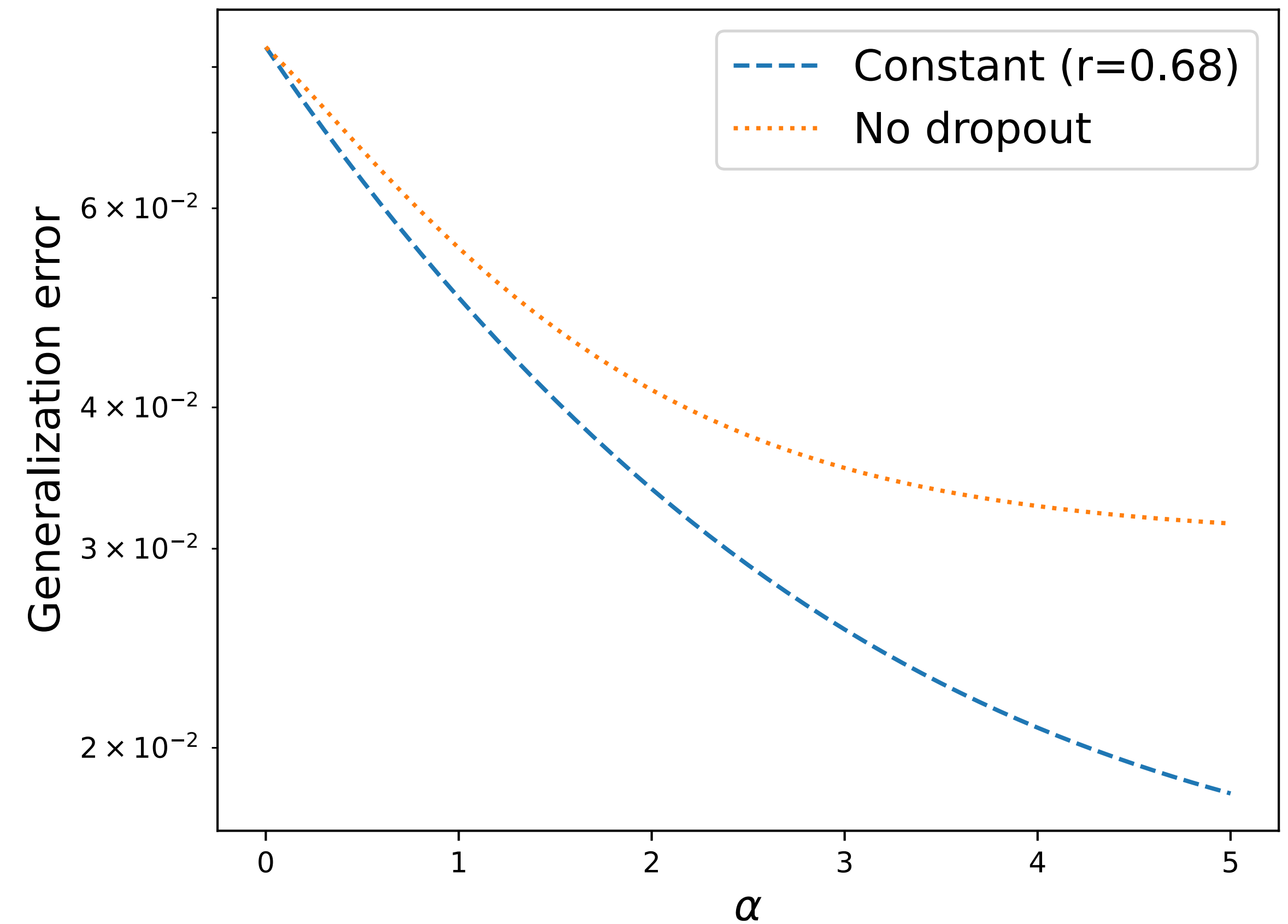
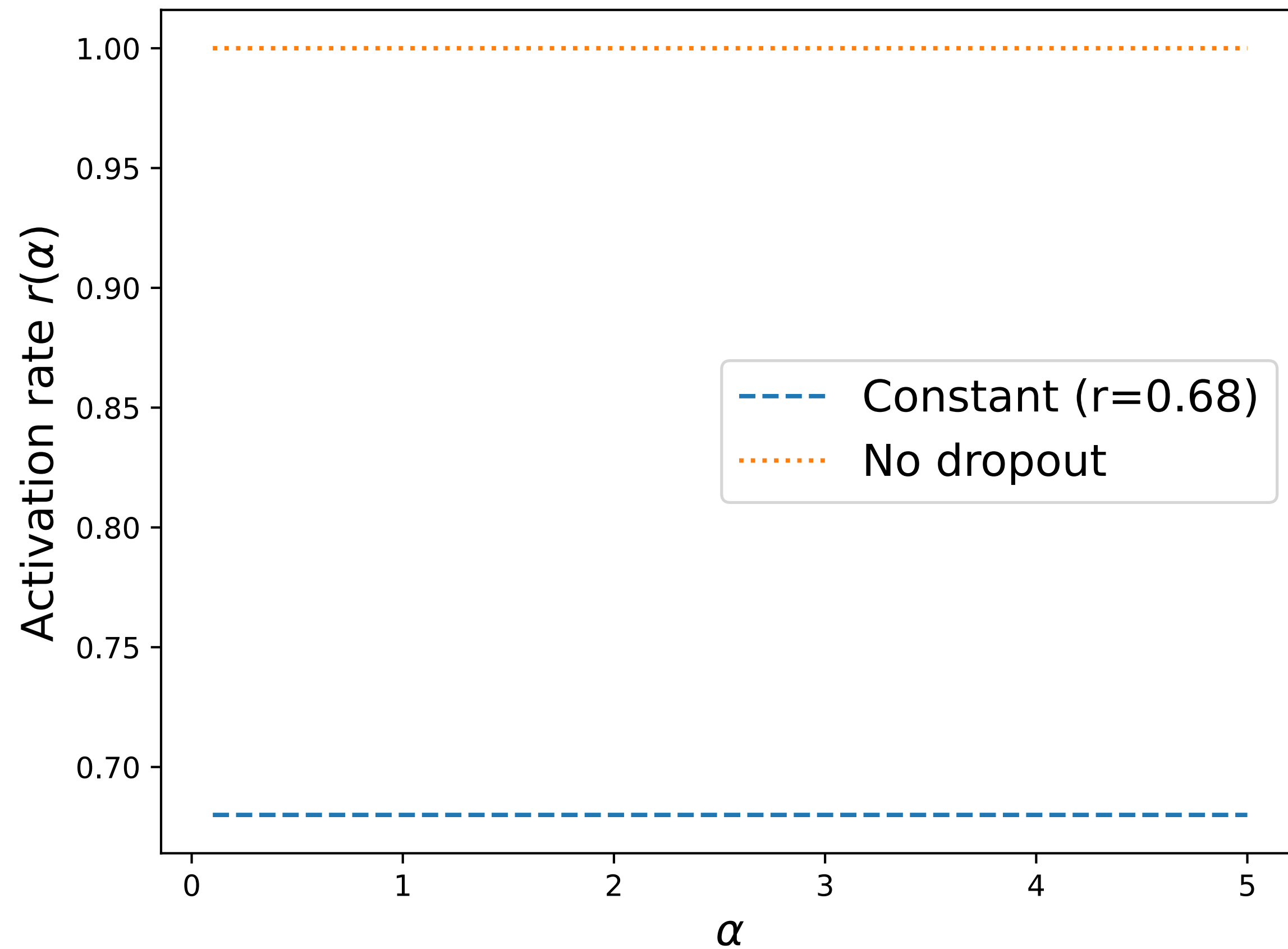
A teacher-student model of dropout

N. Srivastava, G. Hinton, et al., J. ML Res. (2014)

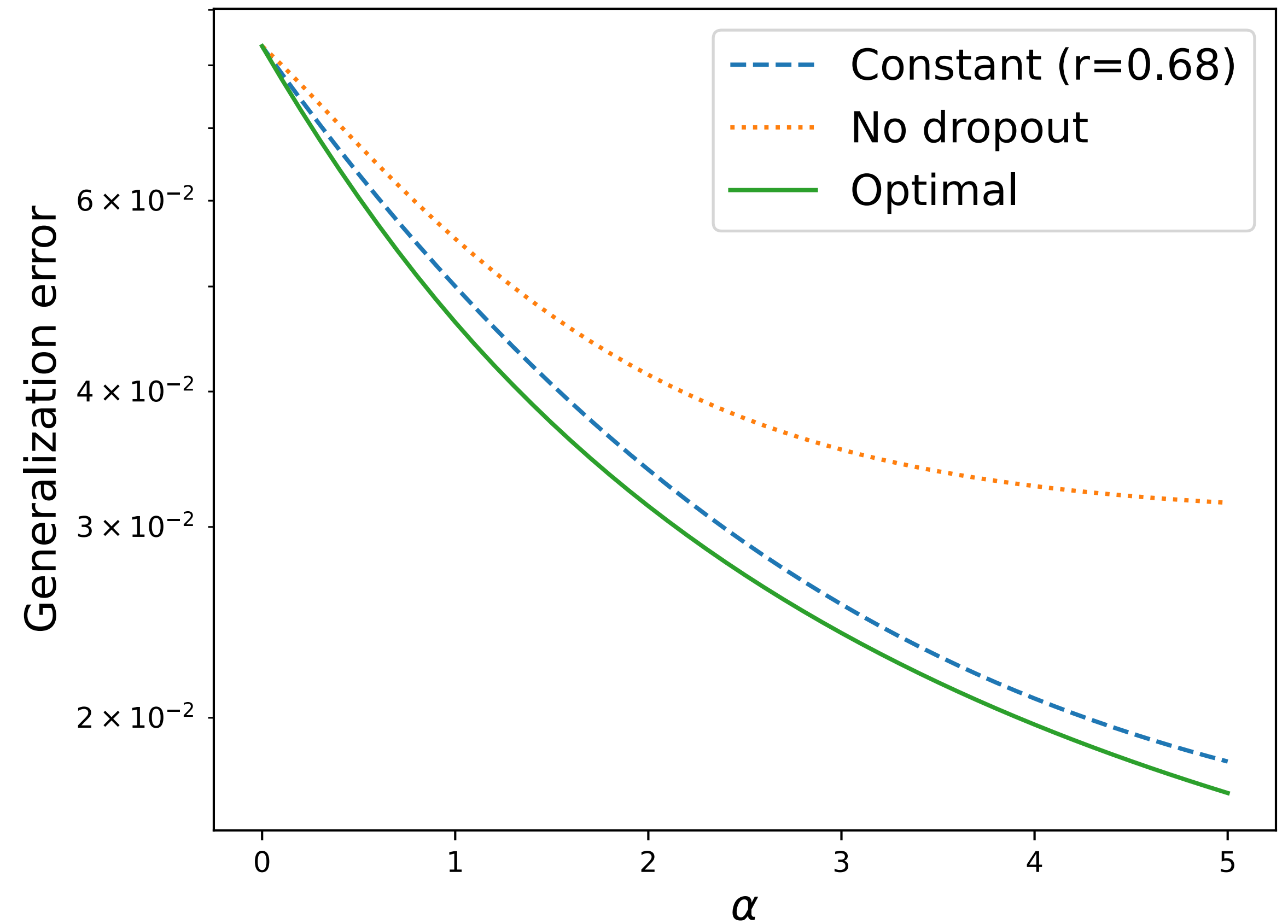
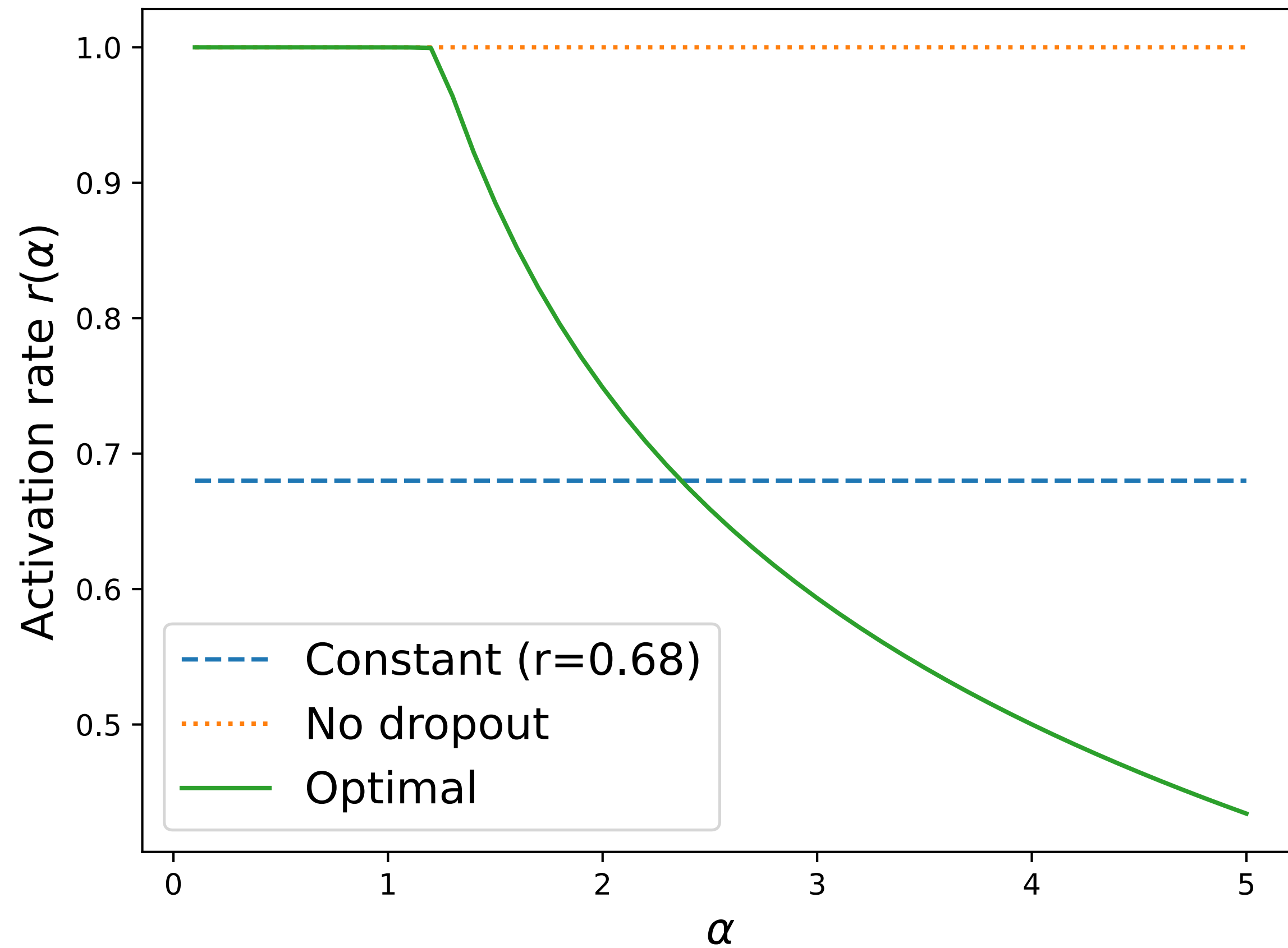


With probability $1 - r$

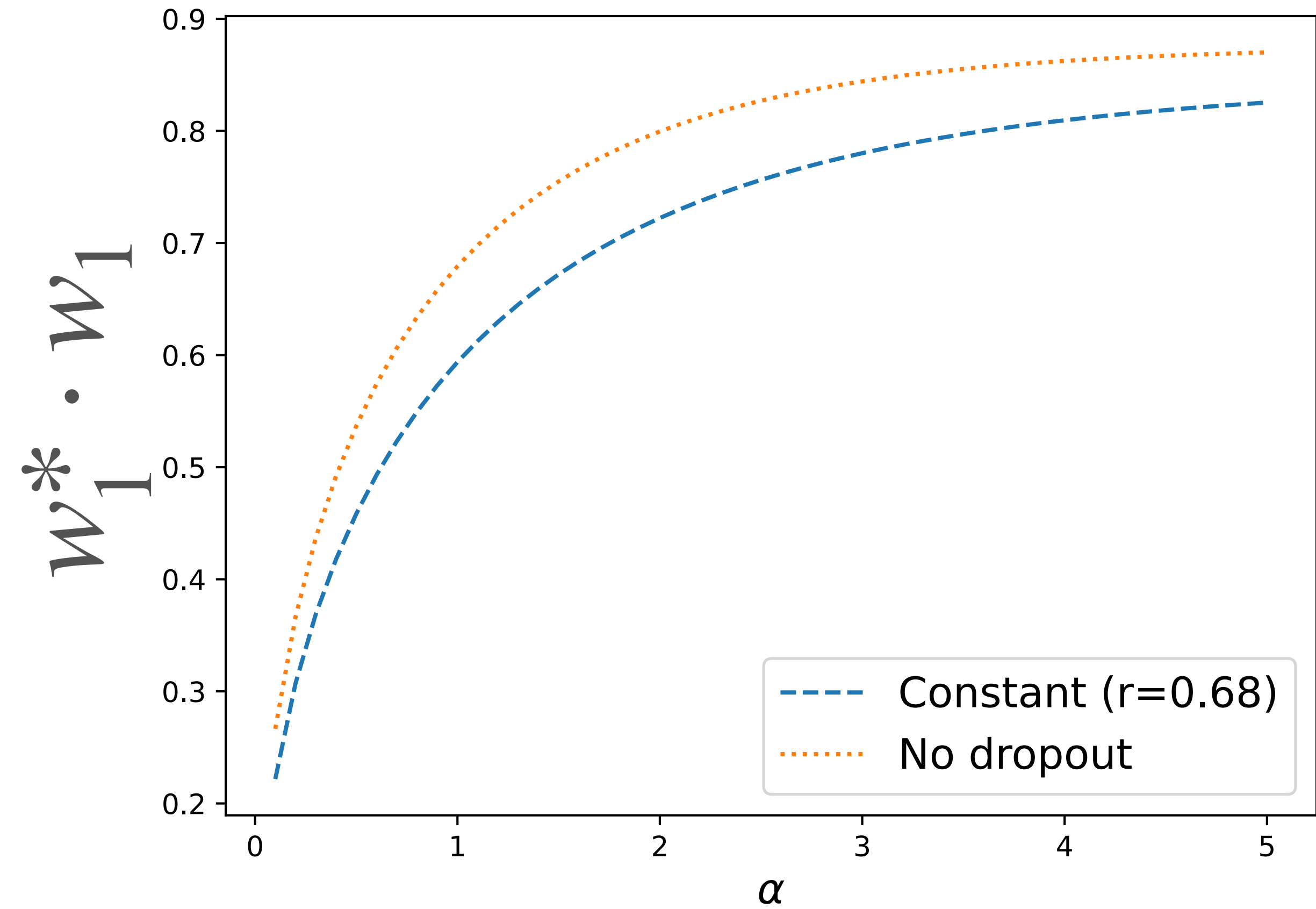
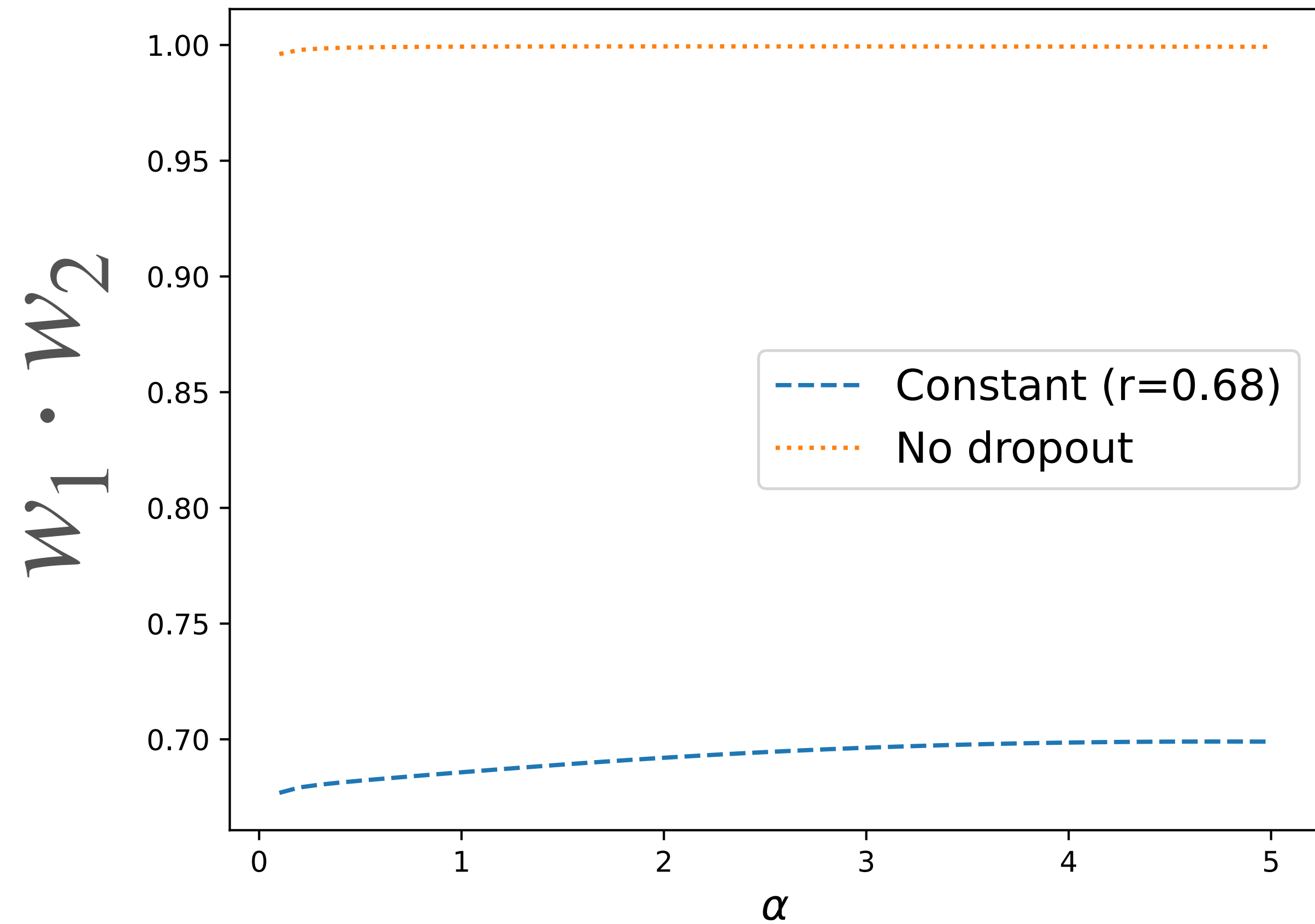
A teacher-student model of dropout



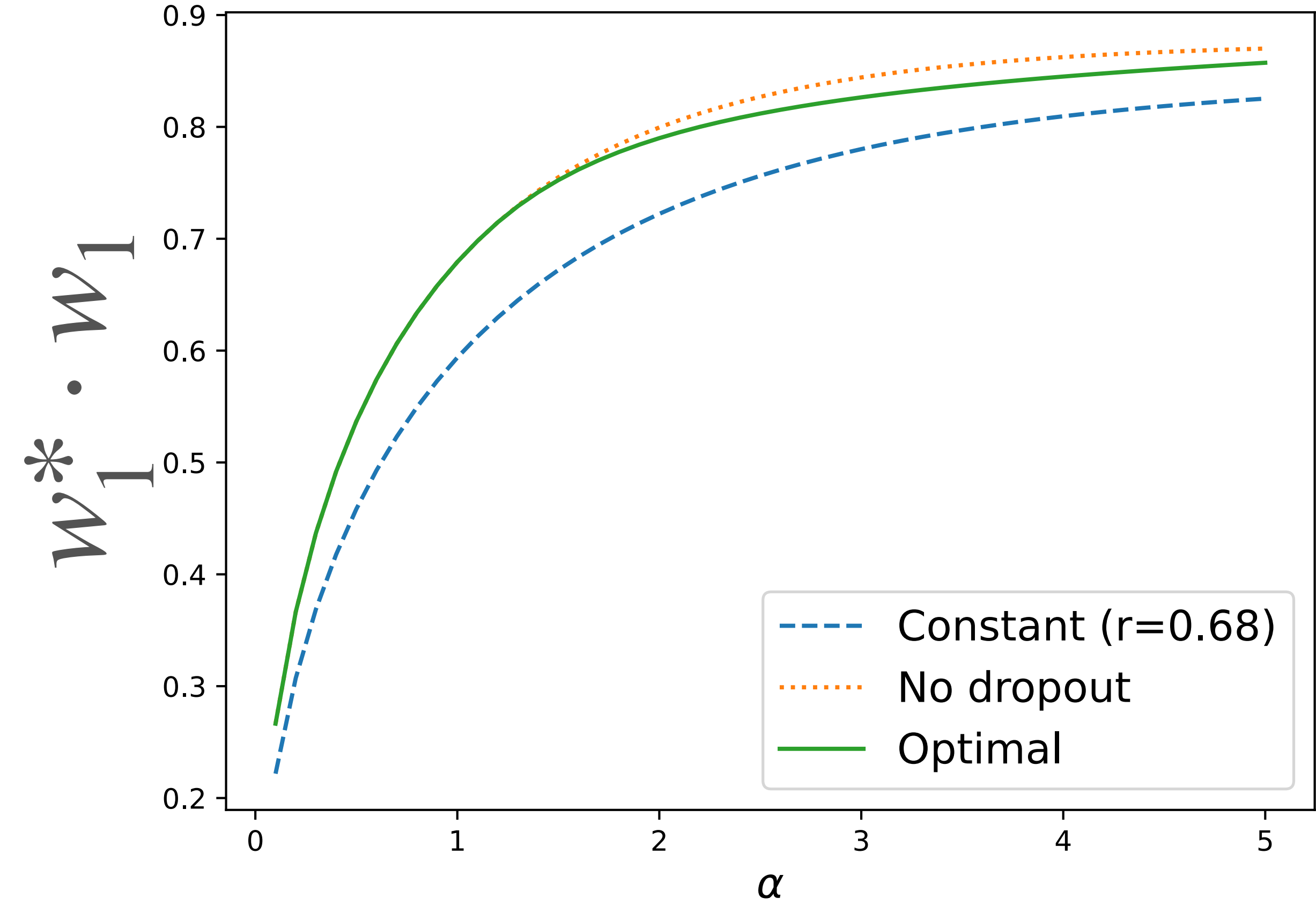
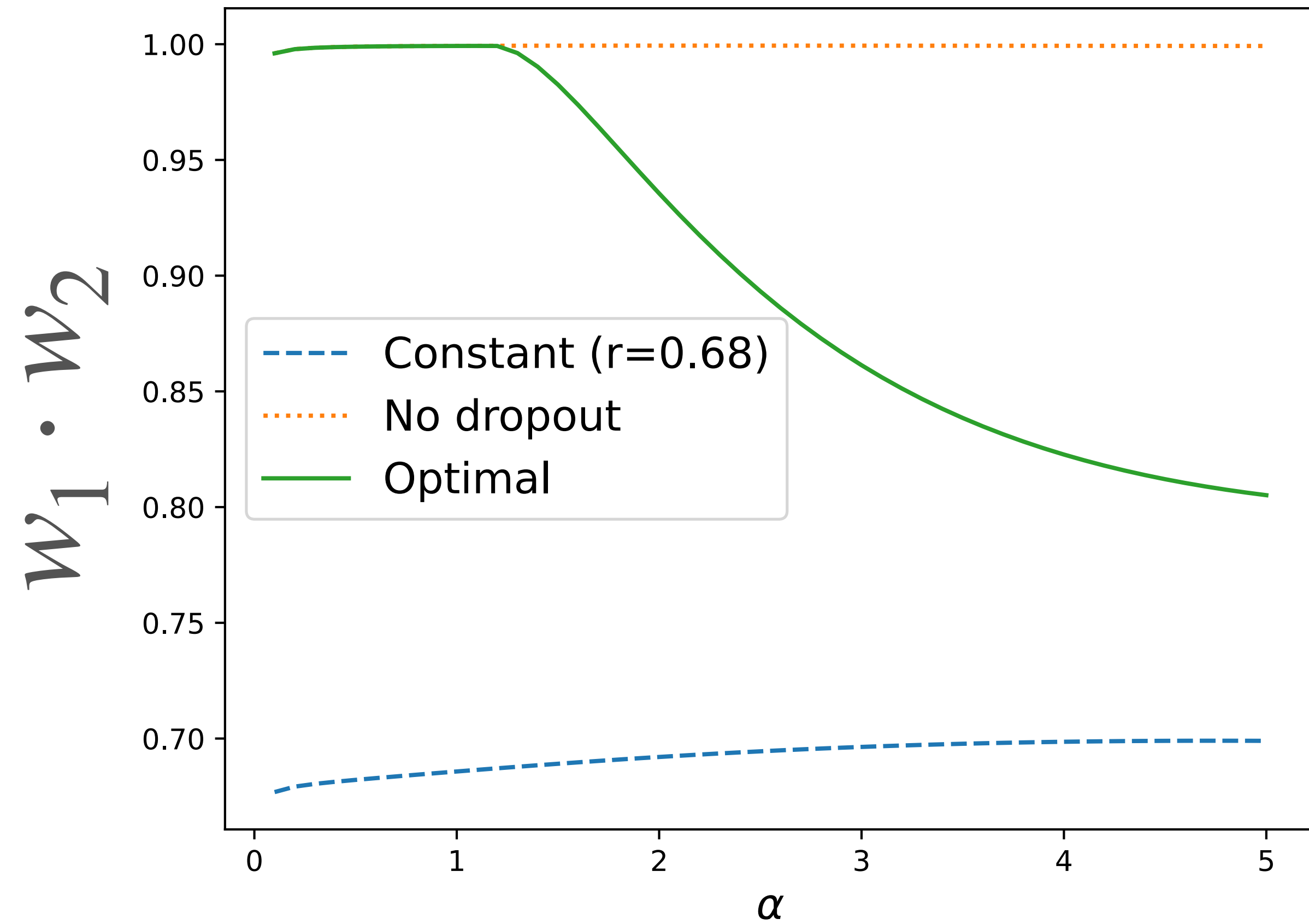
A teacher-student model of dropout



A teacher-student model of dropout



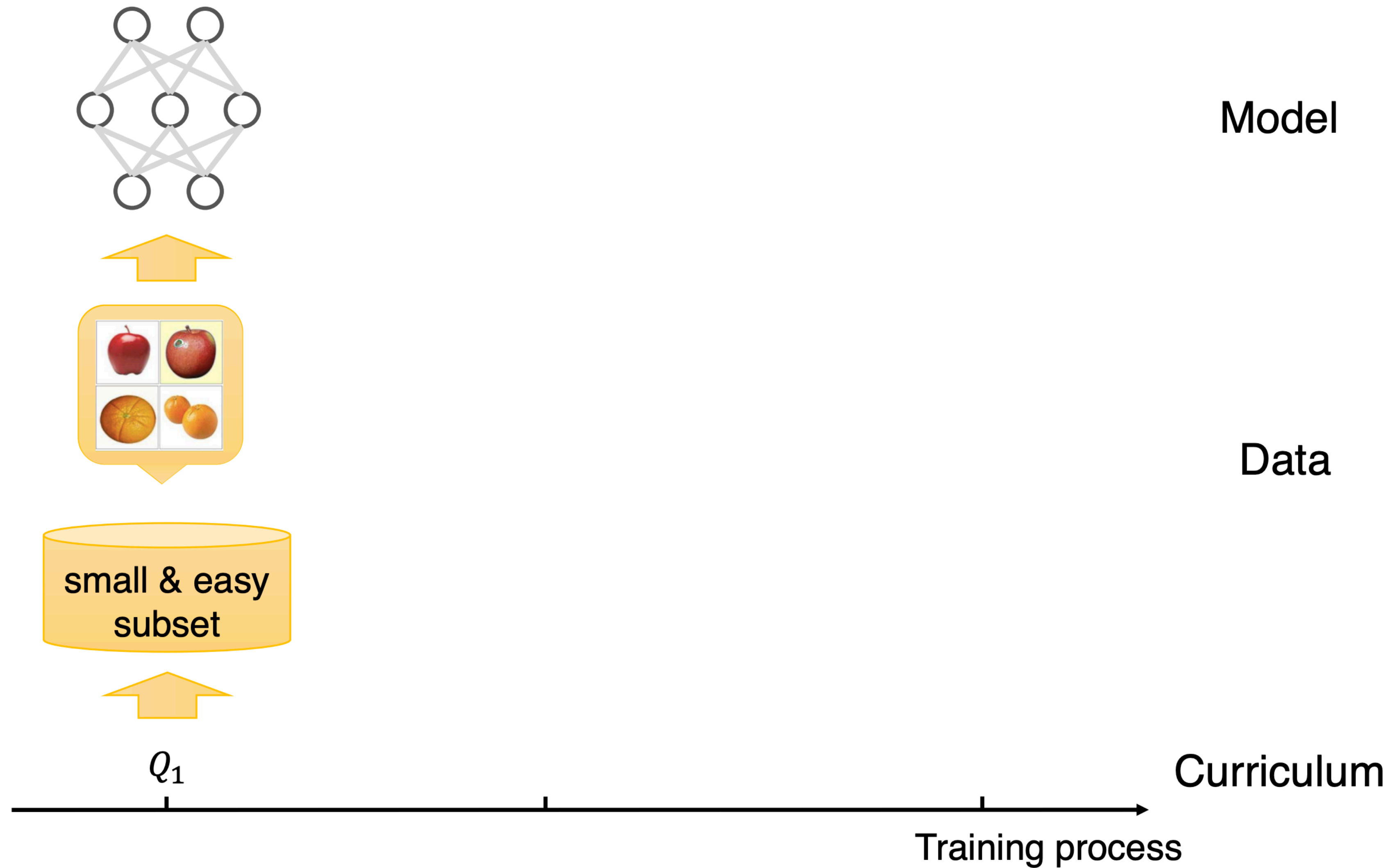
A teacher-student model of dropout



Curriculum learning

(in progress)

Curriculum learning



Model

Data

Curriculum

Training process

Image from: Wang, Xin, Yudong Chen, and Wenwu Zhu. *IEEE transactions on pattern analysis and machine intelligence* 44.9 (2021):4555-4576.

Curriculum learning

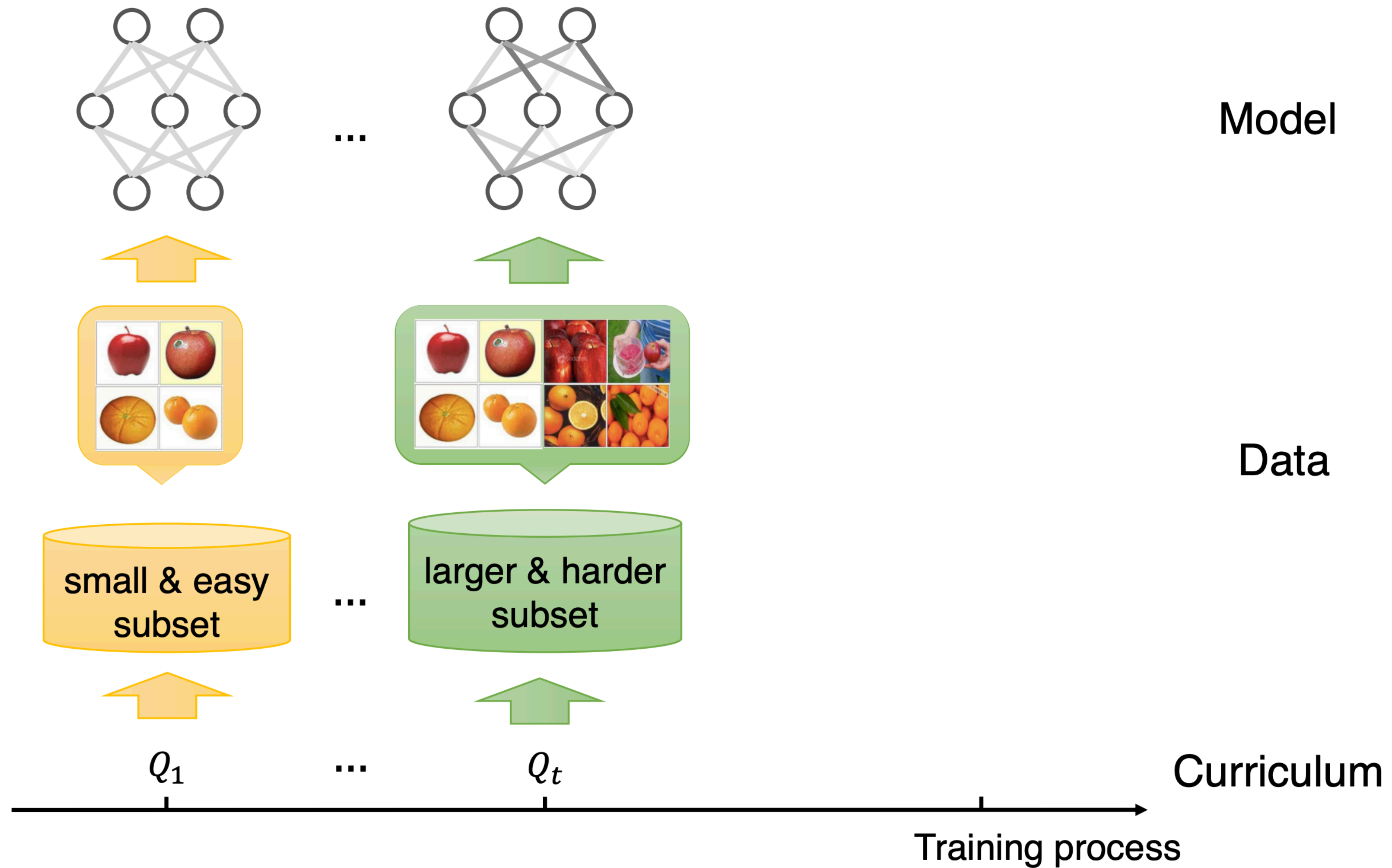


Image from: Wang, Xin, Yudong Chen, and Wenwu Zhu. *IEEE transactions on pattern analysis and machine intelligence* 44.9 (2021):4555-4576.

Curriculum learning

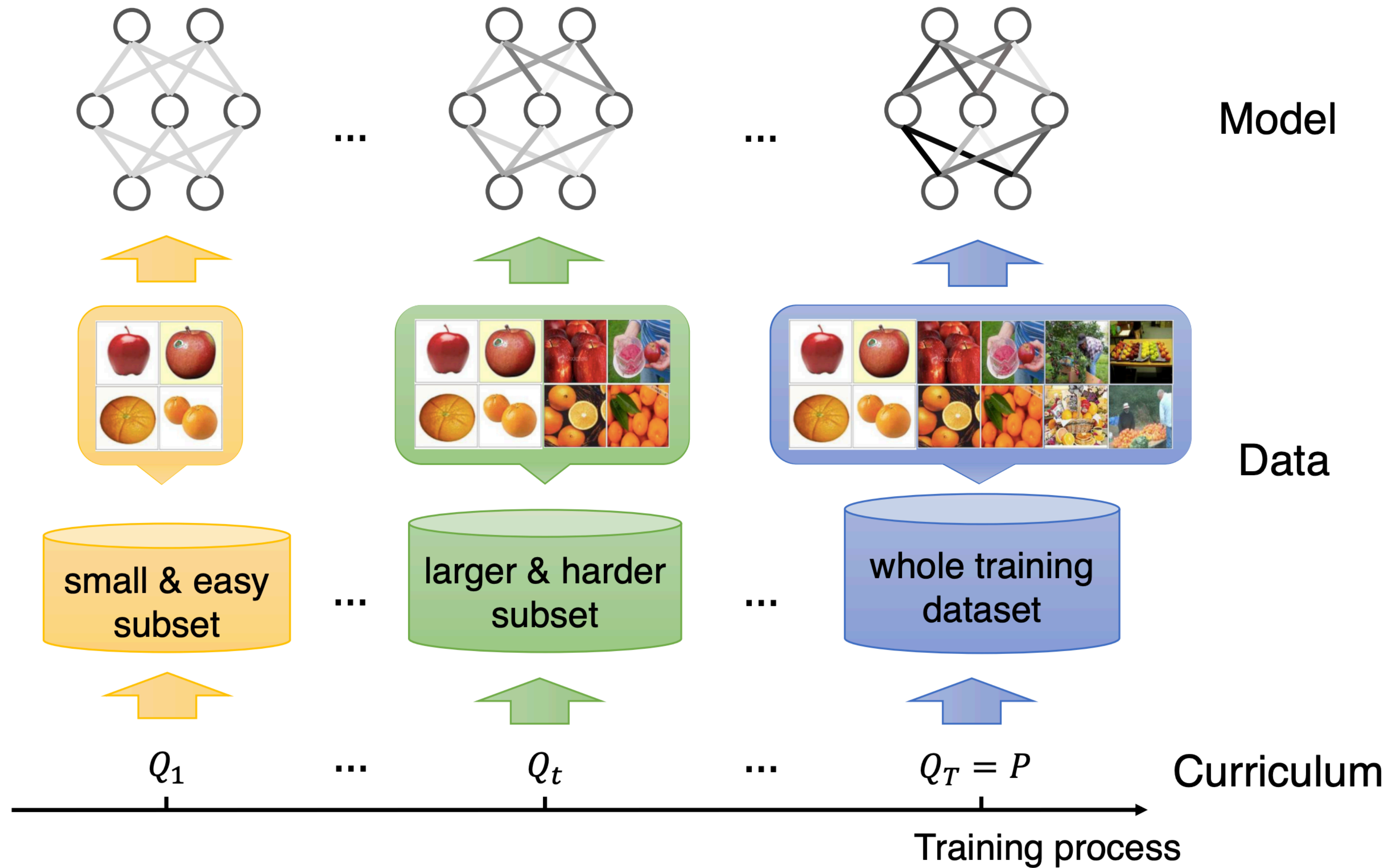


Image from: Wang, Xin, Yudong Chen, and Wenwu Zhu. *IEEE transactions on pattern analysis and machine intelligence* 44.9 (2021):4555-4576.

Curriculum learning

Animals:

Humans:

Curriculum learning

Animals:

Humans:

ML (empirical):

ML (theory):

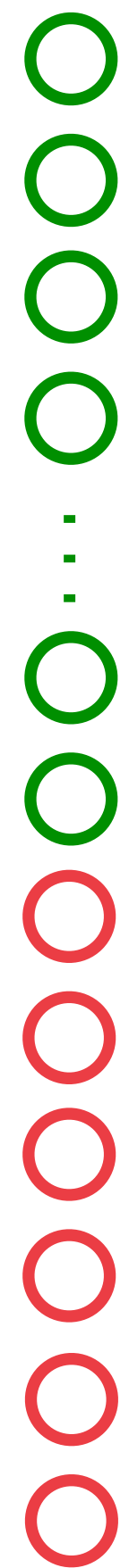
A teacher-student model of curriculum learning

Introduced in: *Bengio, et al. (ICML 2009)*, *Saglietti, et al. (NeurIPS 2022)*

Input: $\mathbf{x} = (\mathbf{x}_r, \mathbf{x}_i) \in \mathbb{R}^N$

Relevant

Irrelevant



$$\mathbf{x}_r \in \mathbb{R}^{\rho N}$$

Unit variance

$$\mathbf{x}_i \in \mathbb{R}^{(1-\rho)N}$$

Variance Δ

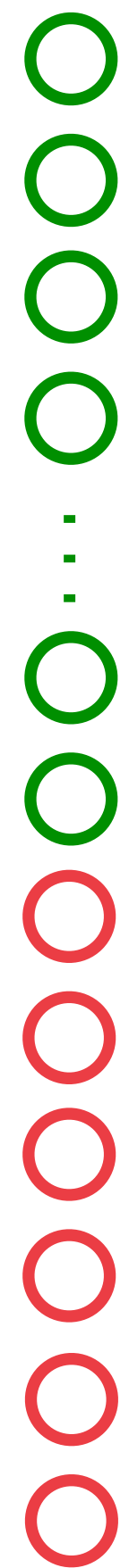
A teacher-student model of curriculum learning

Introduced in: *Bengio, et al. (ICML 2009)*, *Saglietti, et al. (NeurIPS 2022)*

Input: $\mathbf{x} = (\mathbf{x}_r, \mathbf{x}_i) \in \mathbb{R}^N$

Relevant

Irrelevant



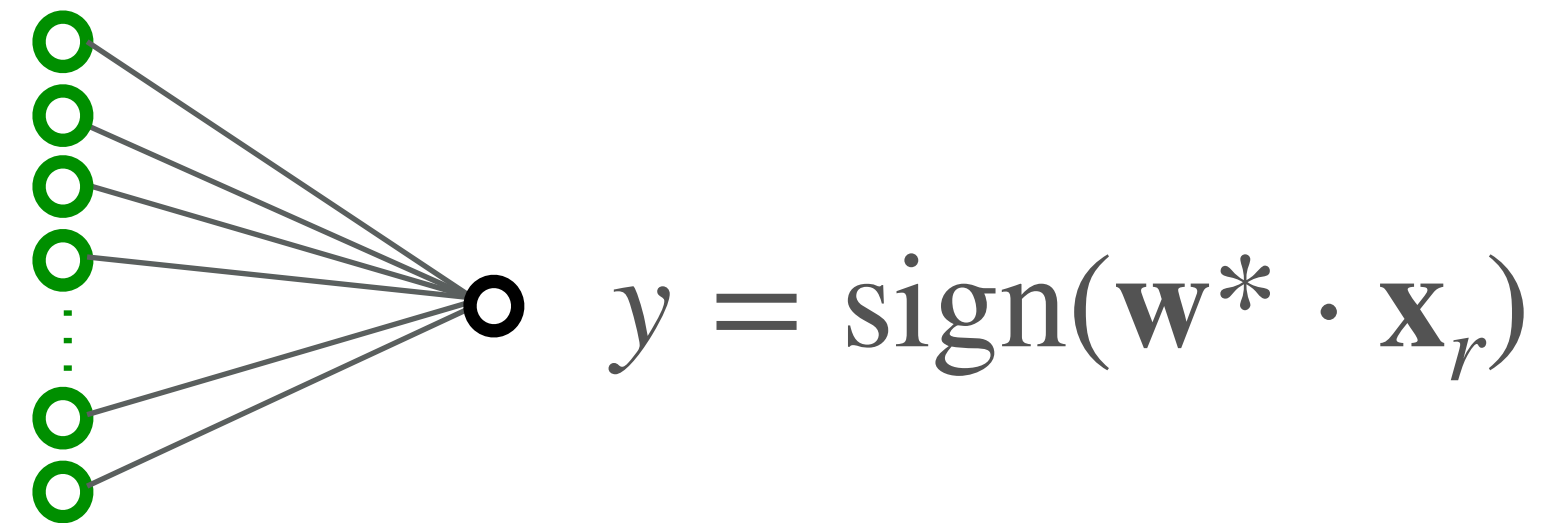
$$\mathbf{x}_r \in \mathbb{R}^{\rho N}$$

Unit variance

$$\mathbf{x}_i \in \mathbb{R}^{(1-\rho)N}$$

Variance Δ

Teacher

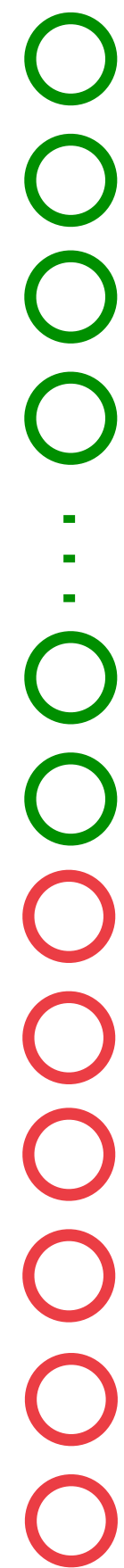


A teacher-student model of curriculum learning

Introduced in: *Bengio, et al. (ICML 2009), Saglietti, et al. (NeurIPS 2022)*

Input: $\mathbf{x} = (\mathbf{x}_r, \mathbf{x}_i) \in \mathbb{R}^N$

Relevant



$\mathbf{x}_r \in \mathbb{R}^{\rho N}$

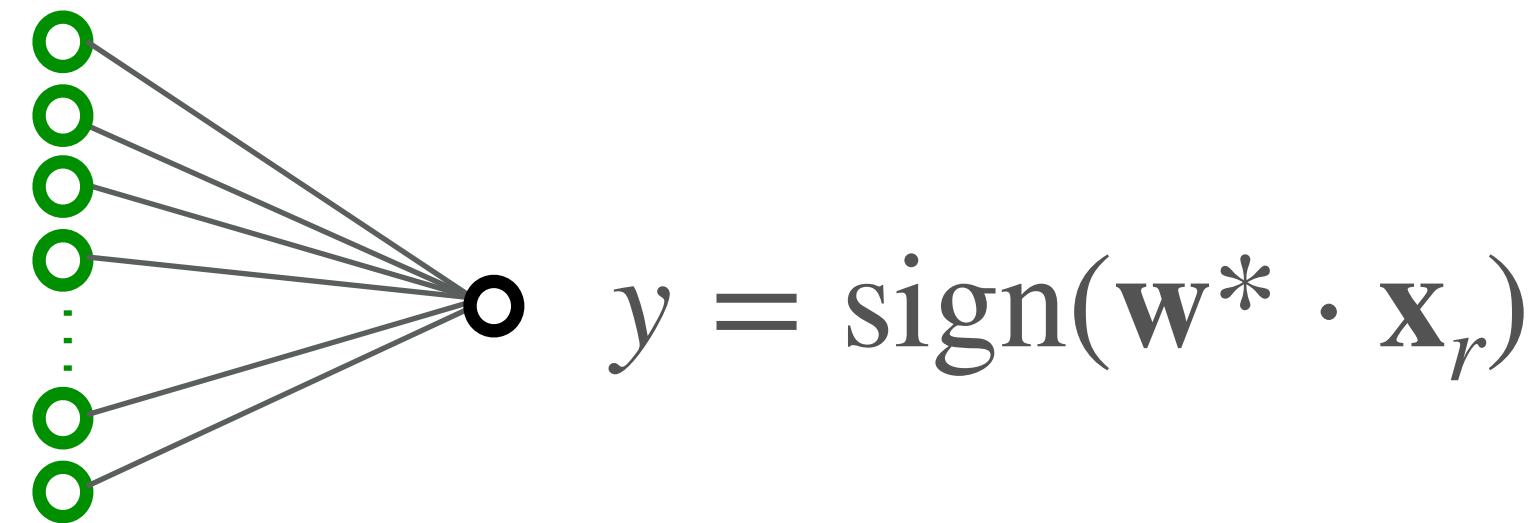
Unit variance

Irrelevant

$\mathbf{x}_i \in \mathbb{R}^{(1-\rho)N}$

Variance Δ

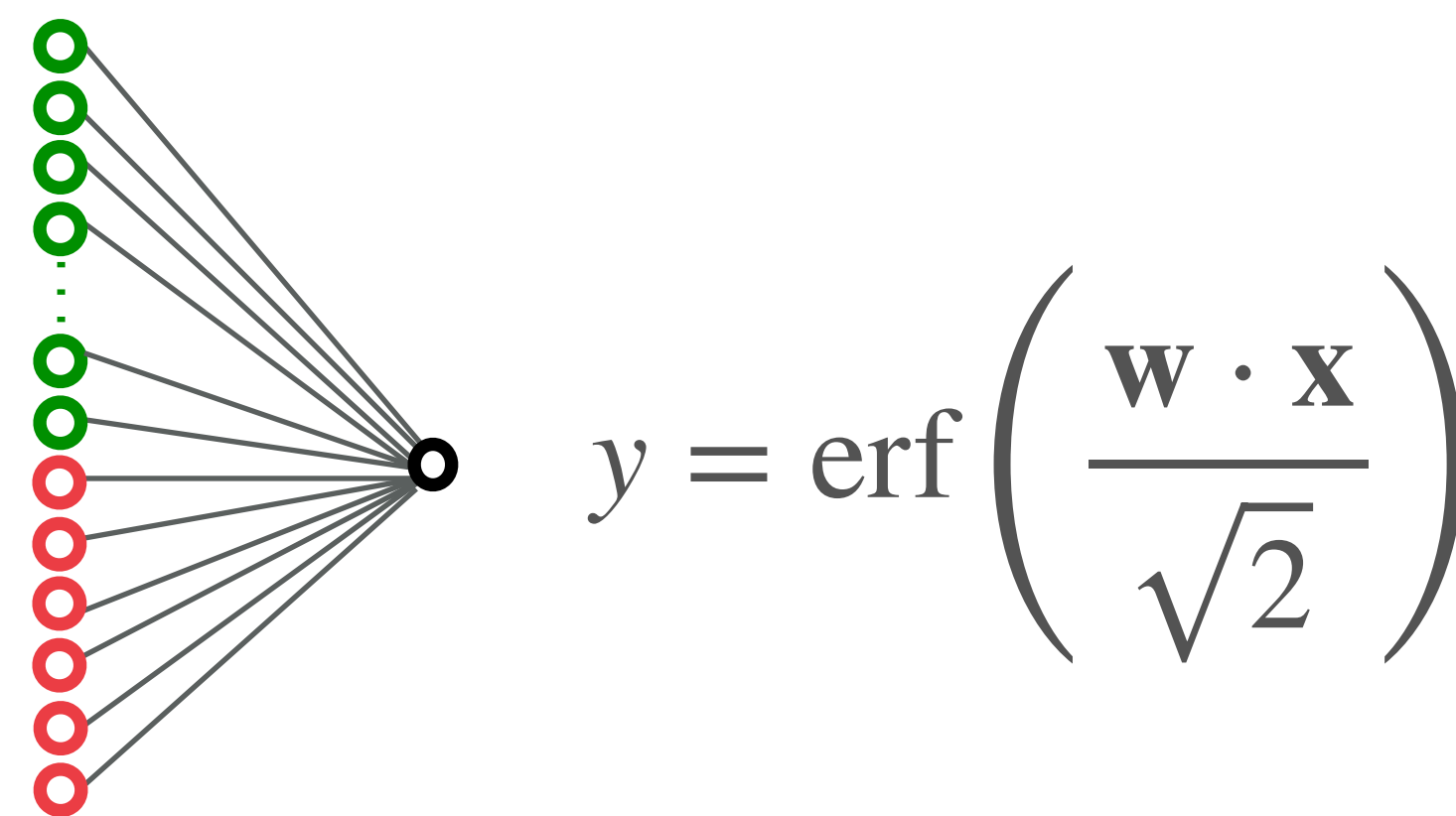
Teacher



Ridge-regularized MSE loss:

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2 + \lambda \|\mathbf{w}\|_2^2$$

Student



An Analytical Theory of Curriculum Learning in Teacher-Student Networks

Luca Saglietti^{†,*}, Stefano Sarao Mannelli^{‡,*}, and Andrew Saxe^{‡,§}

The evolution of the dynamics can be tracked using four order parameters:

$$Q_r = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_r,$$

$$R = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_T,$$

$$Q_i = \frac{1}{N} \mathbf{W}_i \cdot \mathbf{W}_i,$$

$$T = \frac{1}{N} \mathbf{W}_T \cdot \mathbf{W}_T;$$

An Analytical Theory of Curriculum Learning in Teacher-Student Networks

Luca Saglietti^{†,*}, Stefano Sarao Mannelli^{‡,*}, and Andrew Saxe^{‡,§}

The evolution of the dynamics can be tracked using four order parameters:

$$Q_r = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_r,$$

$$R = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_T,$$

$$Q_i = \frac{1}{N} \mathbf{W}_i \cdot \mathbf{W}_i,$$

$$T = \frac{1}{N} \mathbf{W}_T \cdot \mathbf{W}_T;$$

$$Q_r \leftarrow f_{Q_r}(Q_r, Q_i, R, T)$$

$$Q_i \leftarrow f_{Q_i}(Q_r, Q_i, R, T)$$

$$R \leftarrow f_R(Q_r, Q_i, R, T)$$

A teacher-student model of curriculum learning

An Analytical Theory of Curriculum Learning in Teacher-Student Networks

Luca Saglietti^{†,*}, Stefano Sarao Mannelli^{‡,*}, and Andrew Saxe^{‡,§}

The evolution of the dynamics can be tracked using four order parameters:

$$Q_r = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_r,$$

$$R = \frac{1}{N} \mathbf{W}_r \cdot \mathbf{W}_T,$$

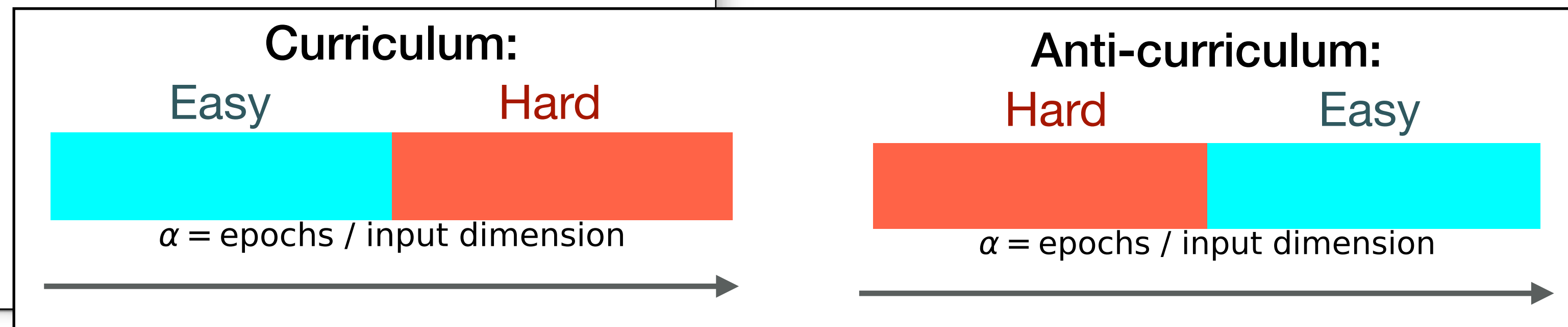
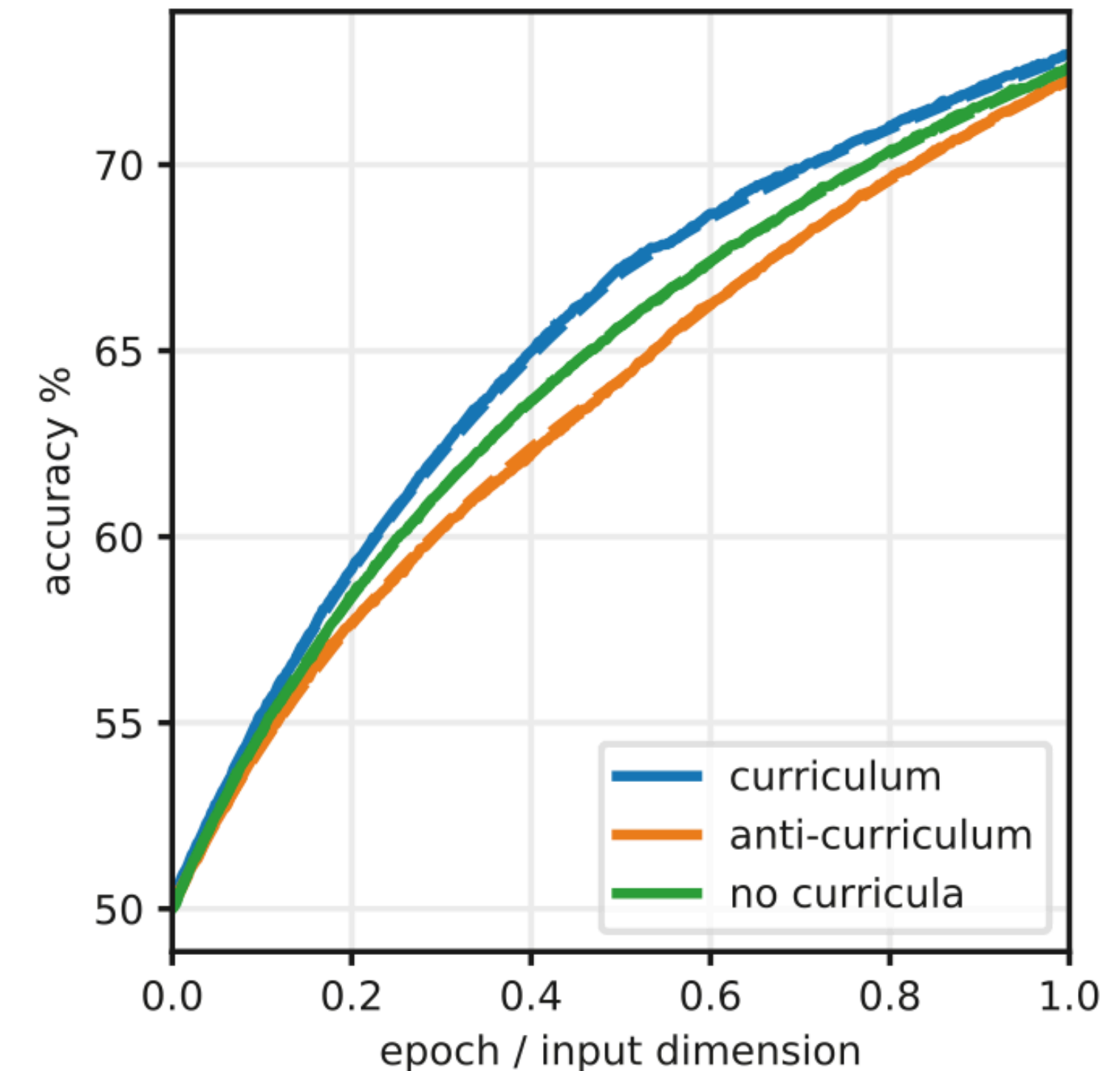
$$Q_i = \frac{1}{N} \mathbf{W}_i \cdot \mathbf{W}_i,$$

$$T = \frac{1}{N} \mathbf{W}_T \cdot \mathbf{W}_T;$$

$$Q_r \leftarrow f_{Q_r}(Q_r, Q_i, R, T)$$

$$Q_i \leftarrow f_{Q_i}(Q_r, Q_i, R, T)$$

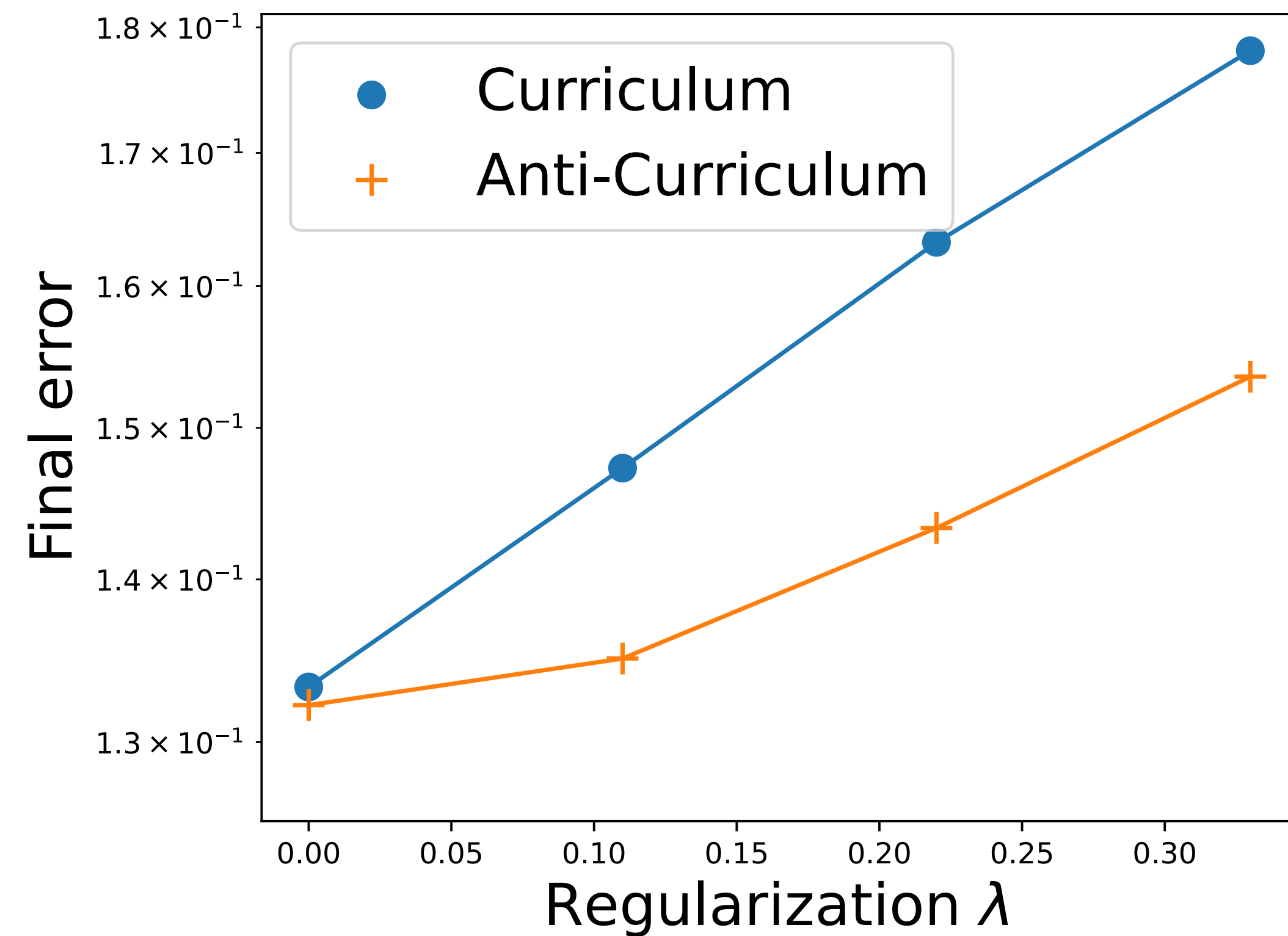
$$R \leftarrow f_R(Q_r, Q_i, R, T)$$



Optimal curriculum protocol

$$\rho = 0.55, \eta = 2.58$$

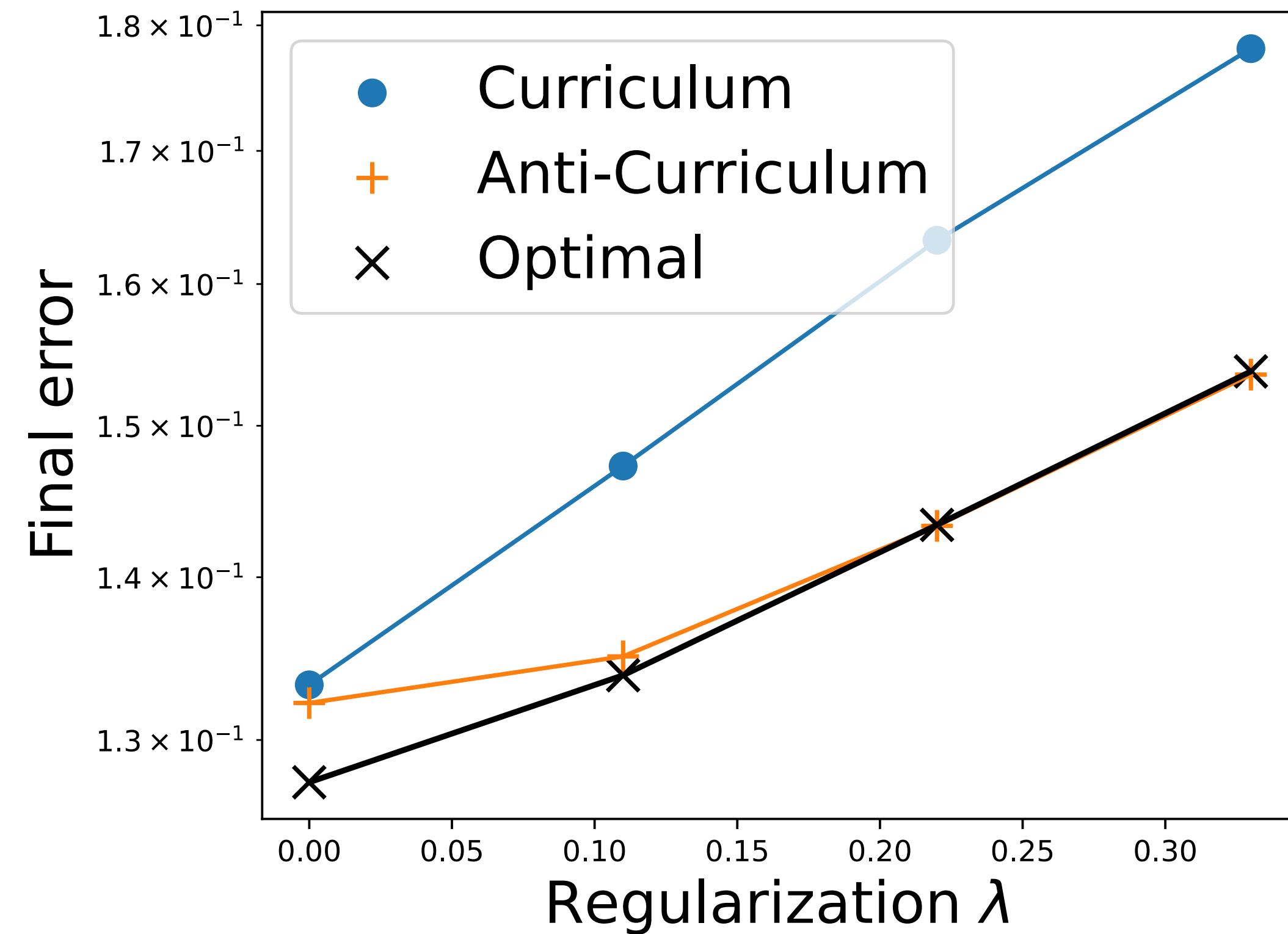
Control: $\mathbf{u} = \Delta$



Optimal curriculum protocol

$$\rho = 0.55, \eta = 2.58$$

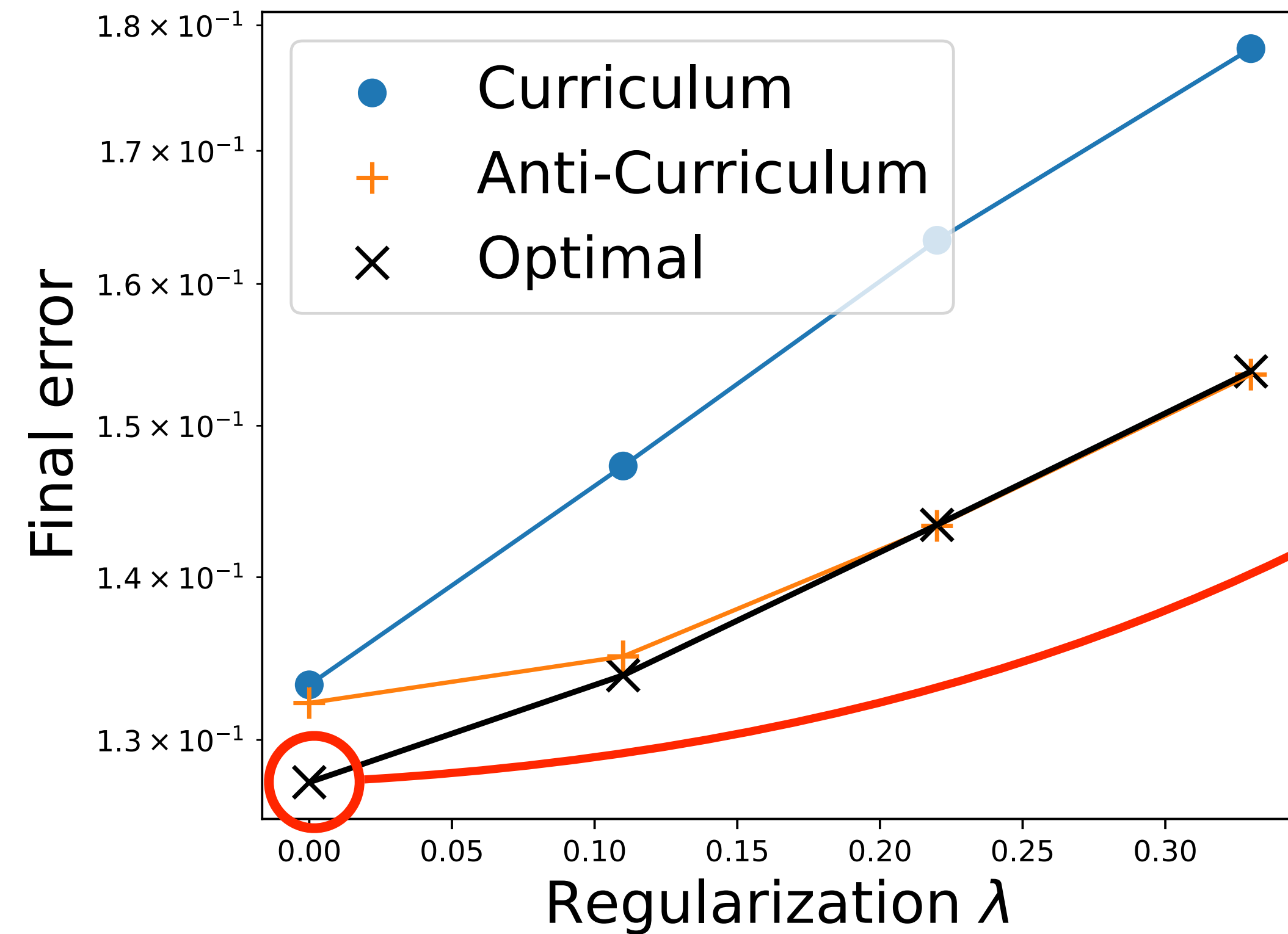
Control: $\mathbf{u} = \Delta$



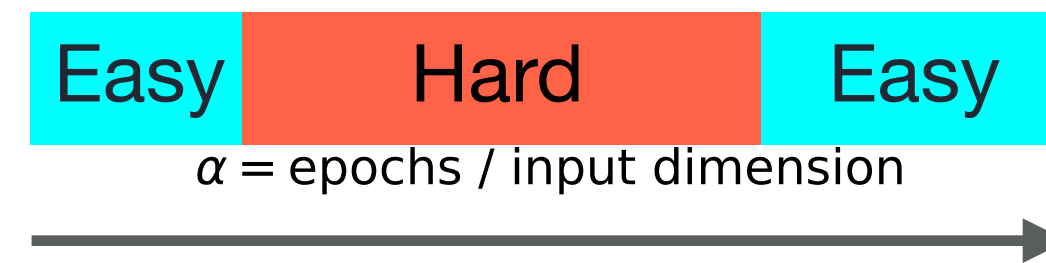
Optimal curriculum protocol

$$\rho = 0.55, \eta = 2.58$$

Control: $\mathbf{u} = \Delta$



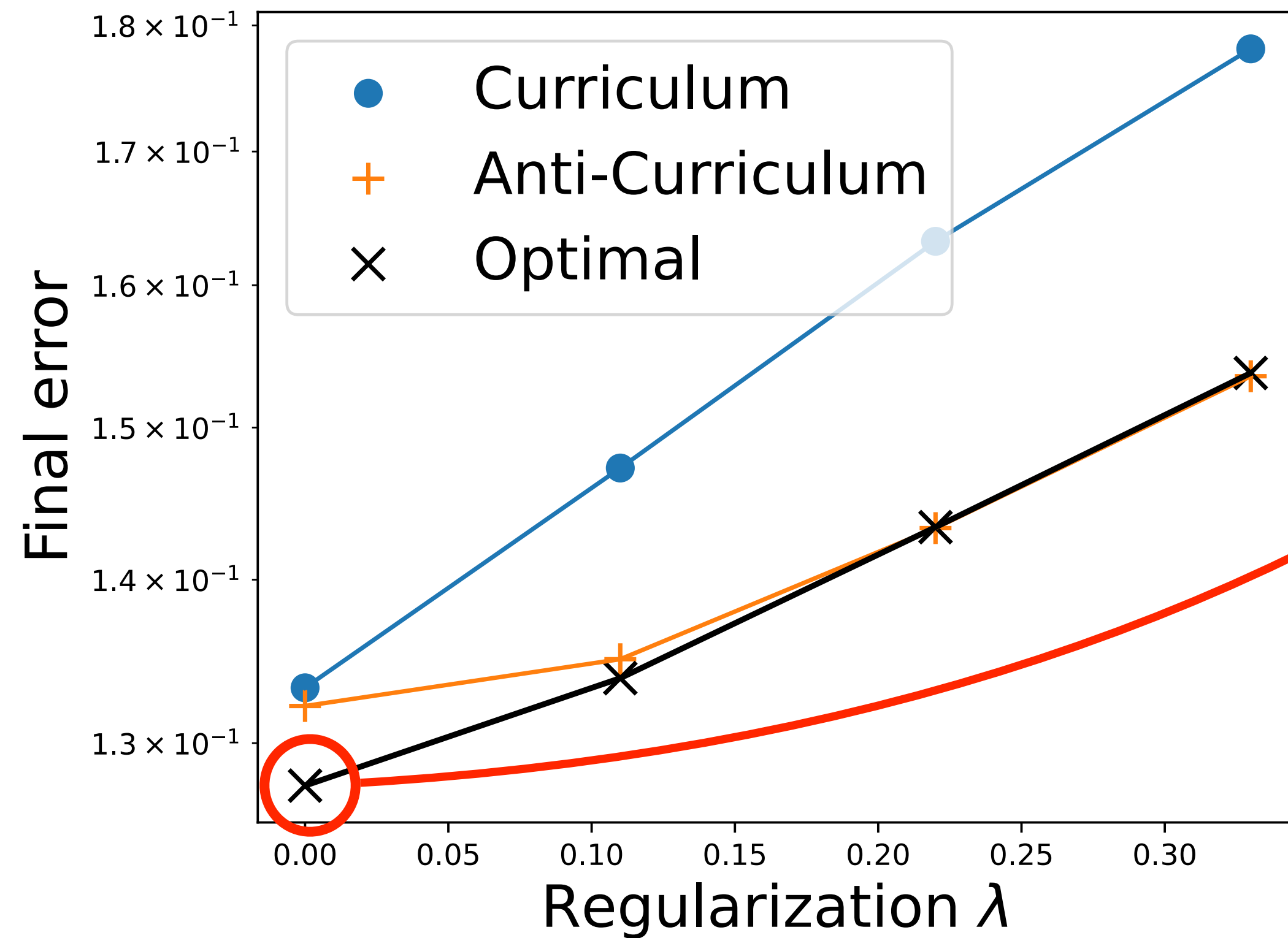
Non-monotonic
curriculum is optimal:



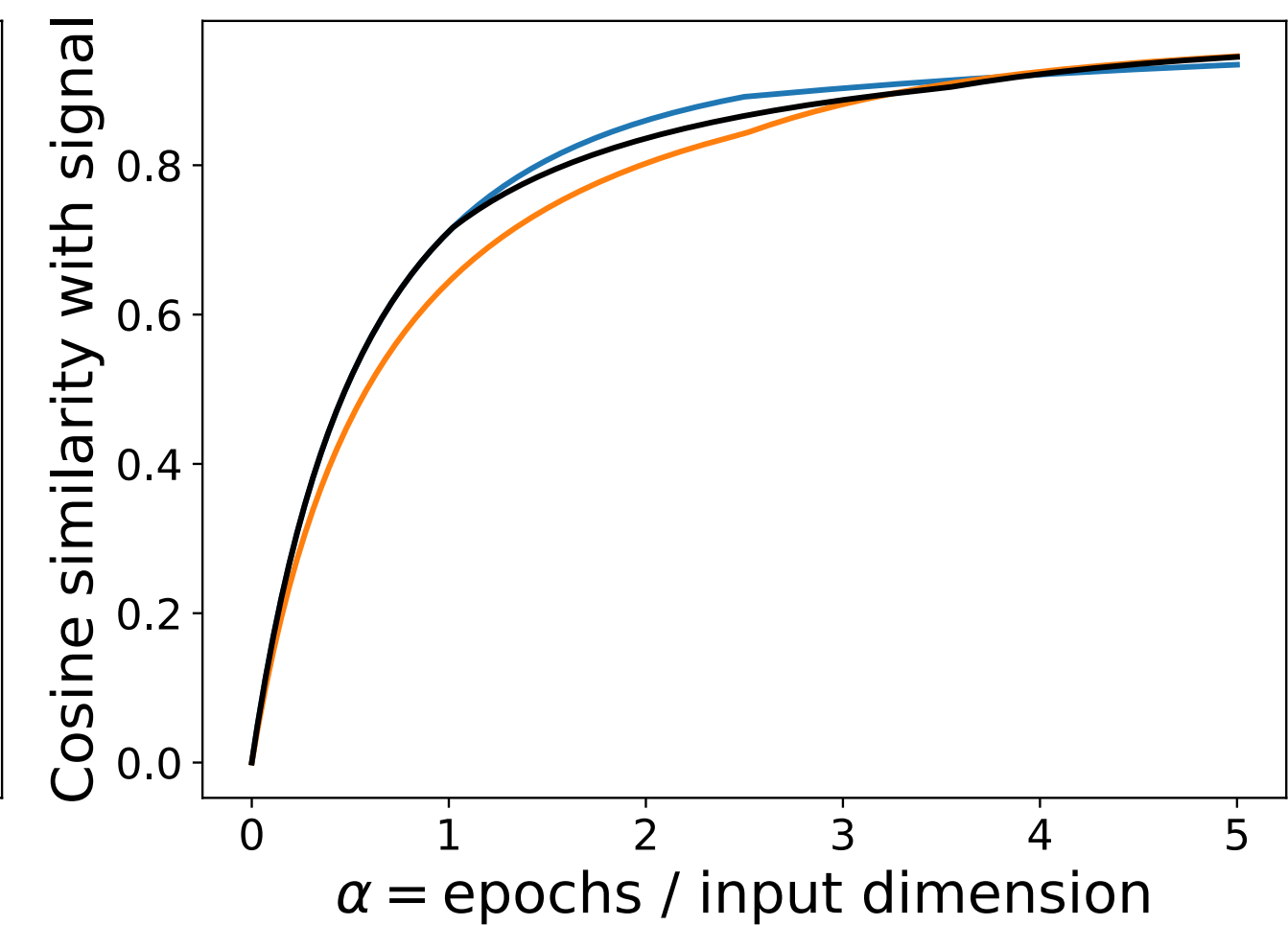
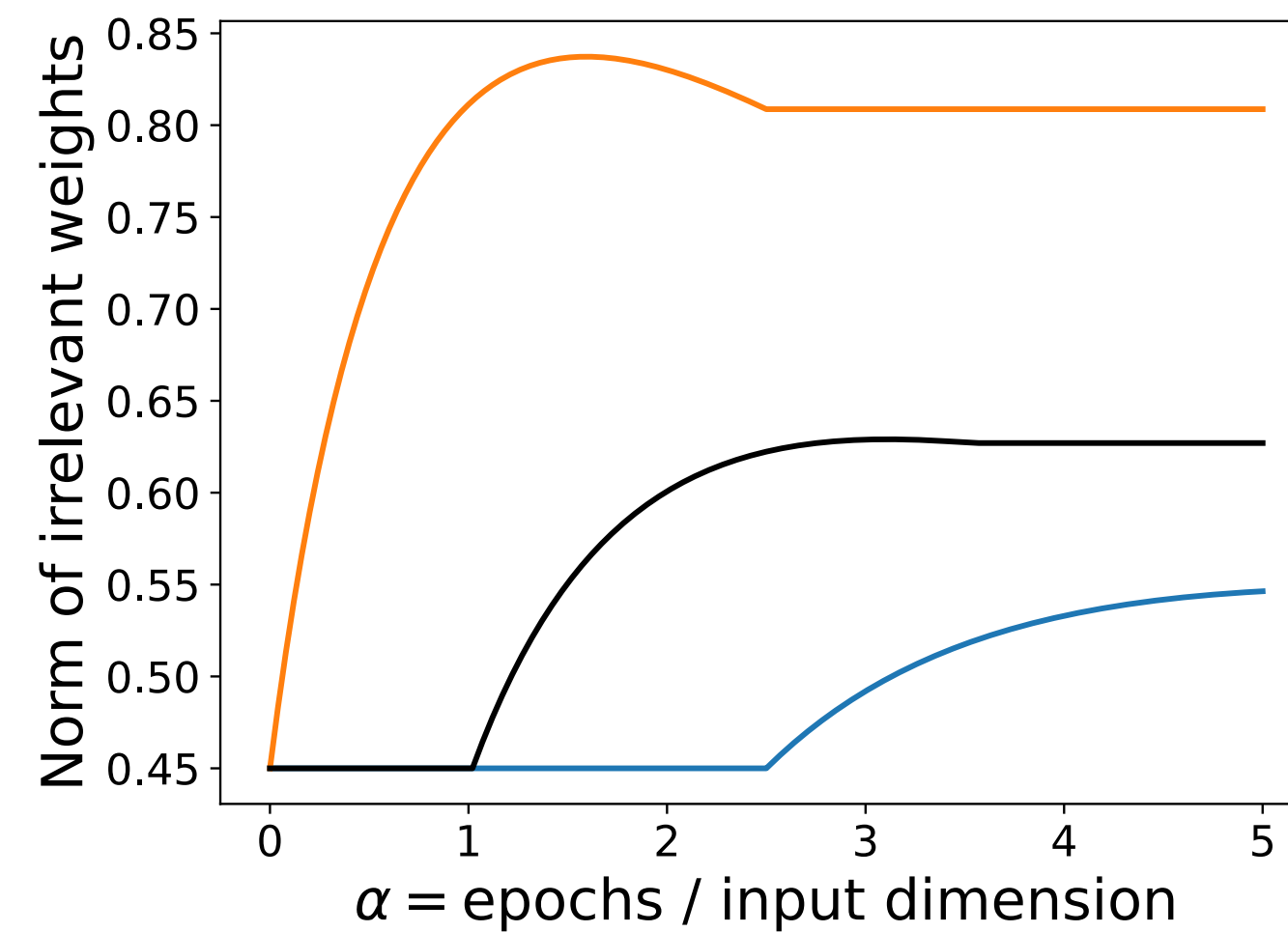
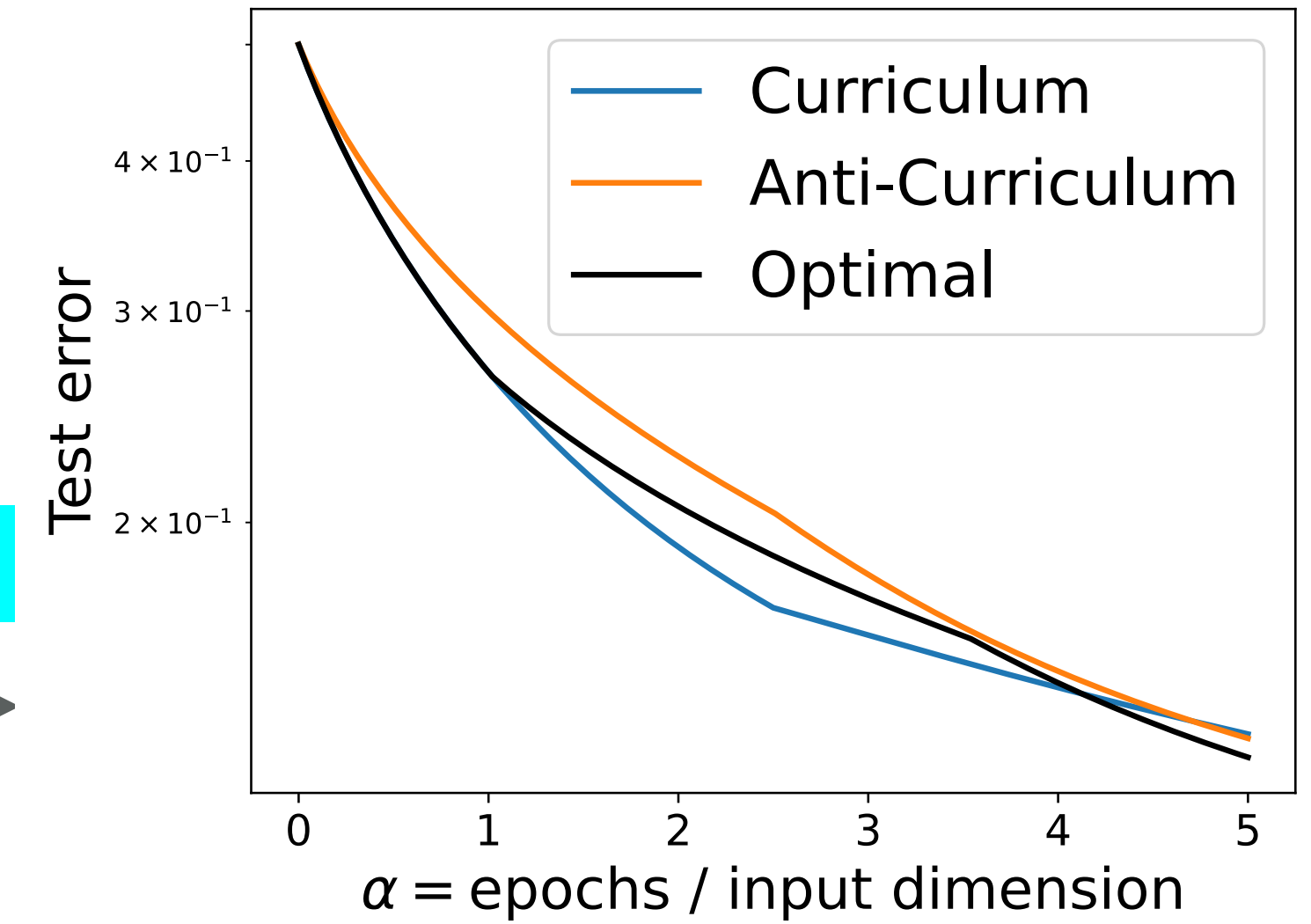
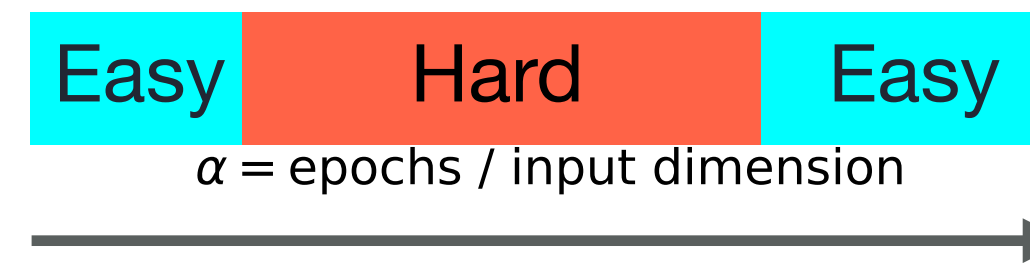
Optimal curriculum protocol

$$\rho = 0.55, \eta = 2.58$$

Control: $\mathbf{u} = \Delta$



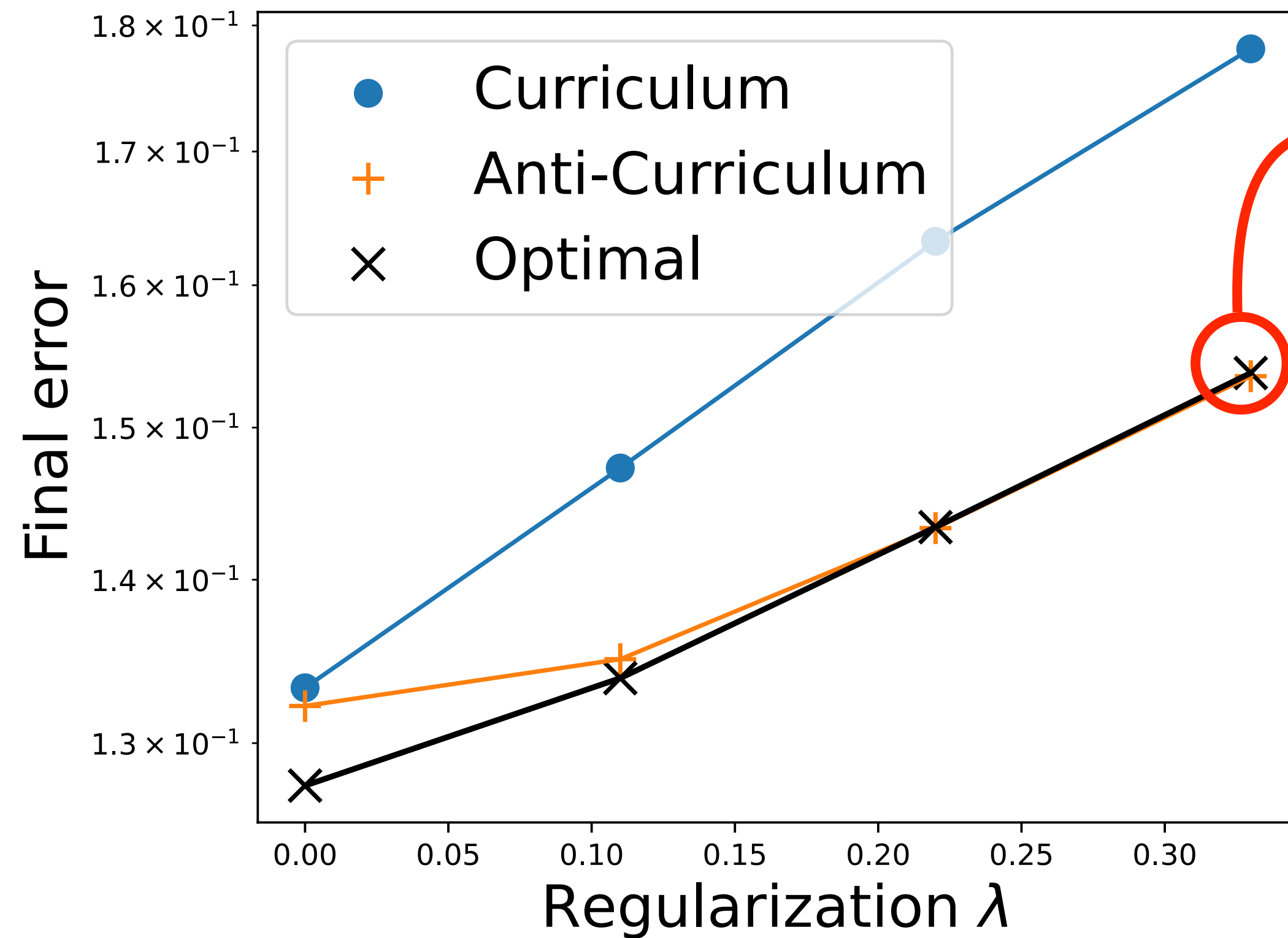
Non-monotonic
curriculum is optimal:



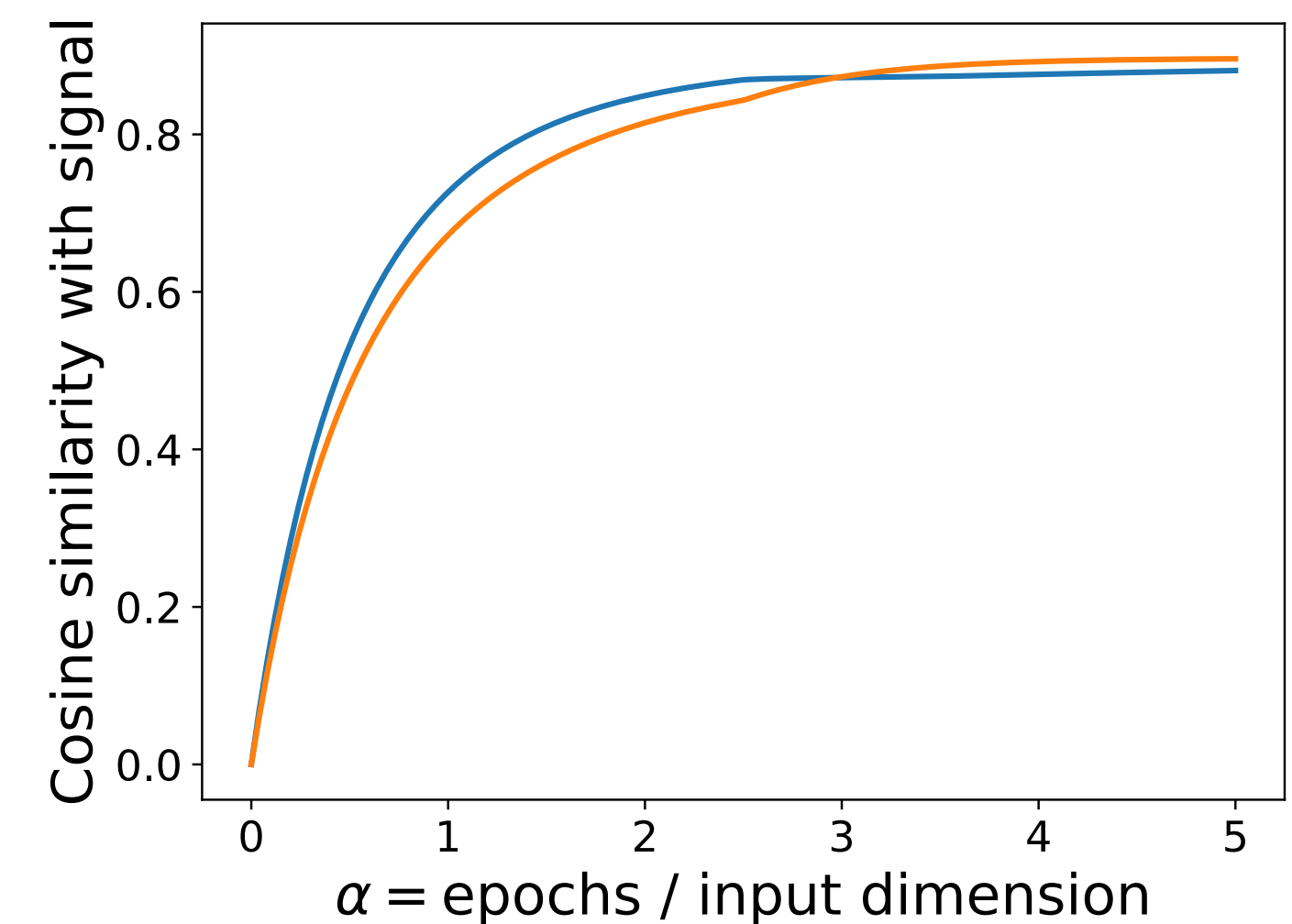
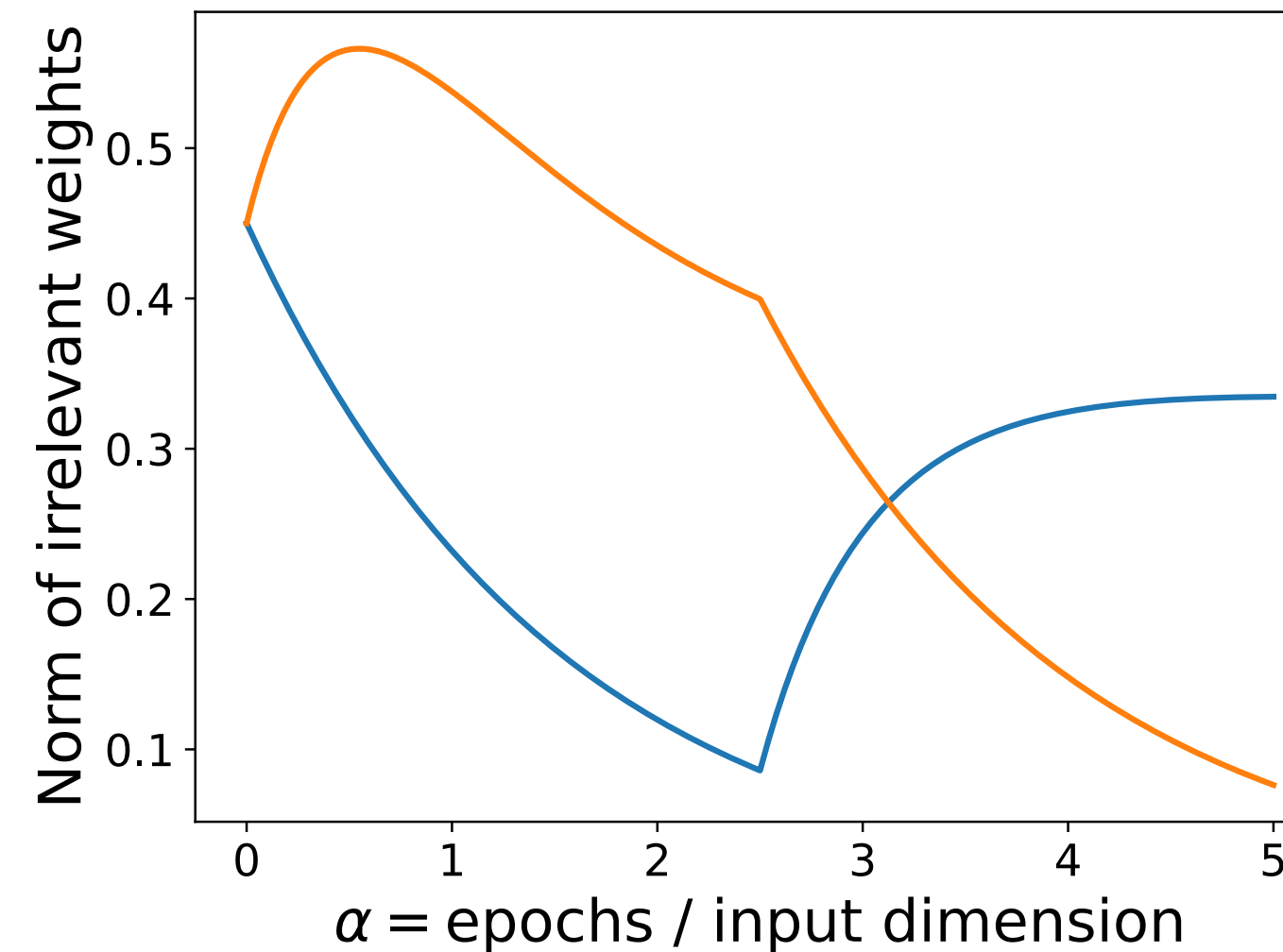
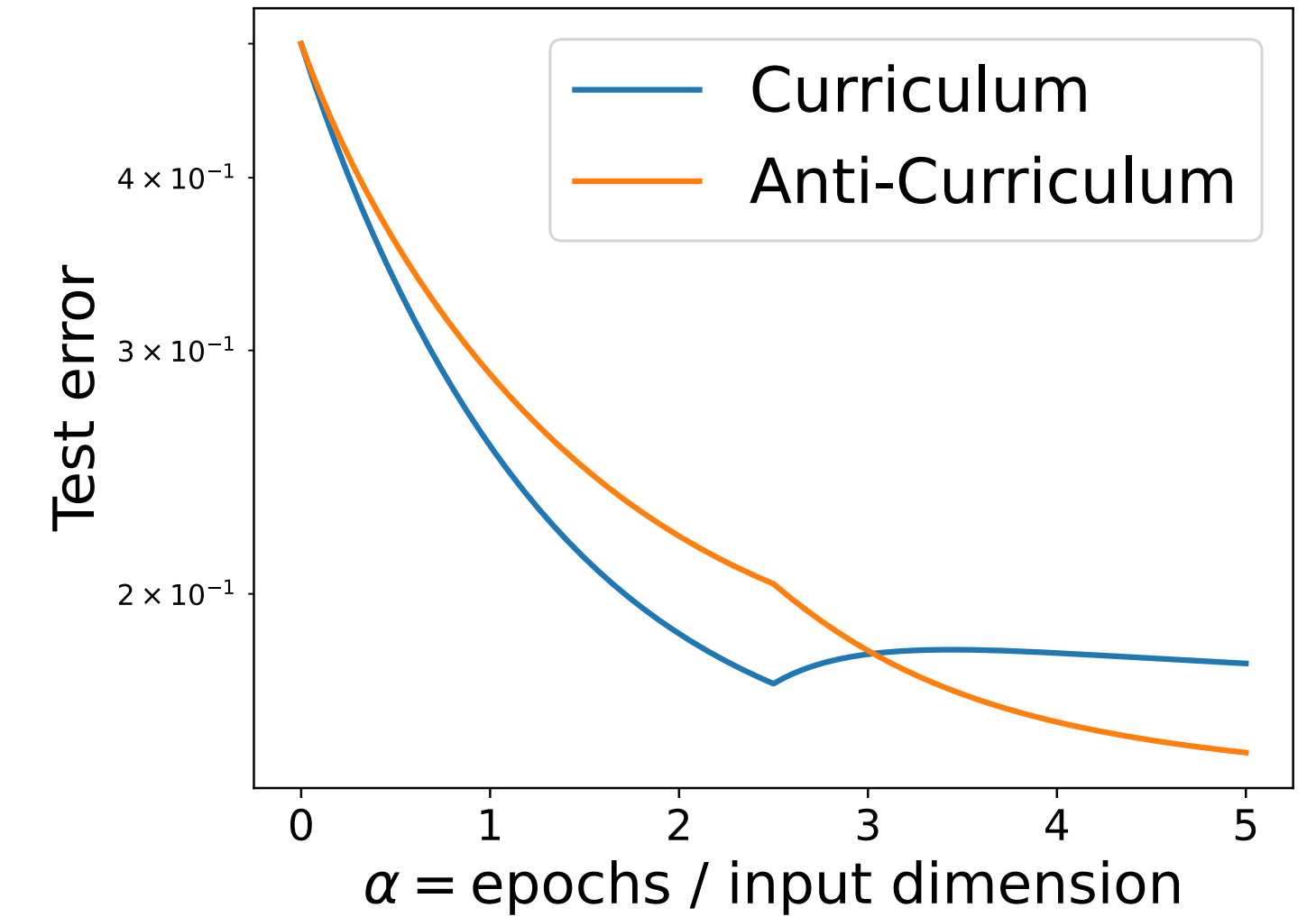
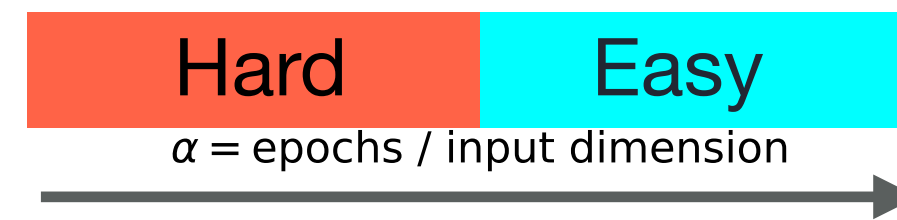
Optimal curriculum protocol

$$\rho = 0.55, \eta = 2.58$$

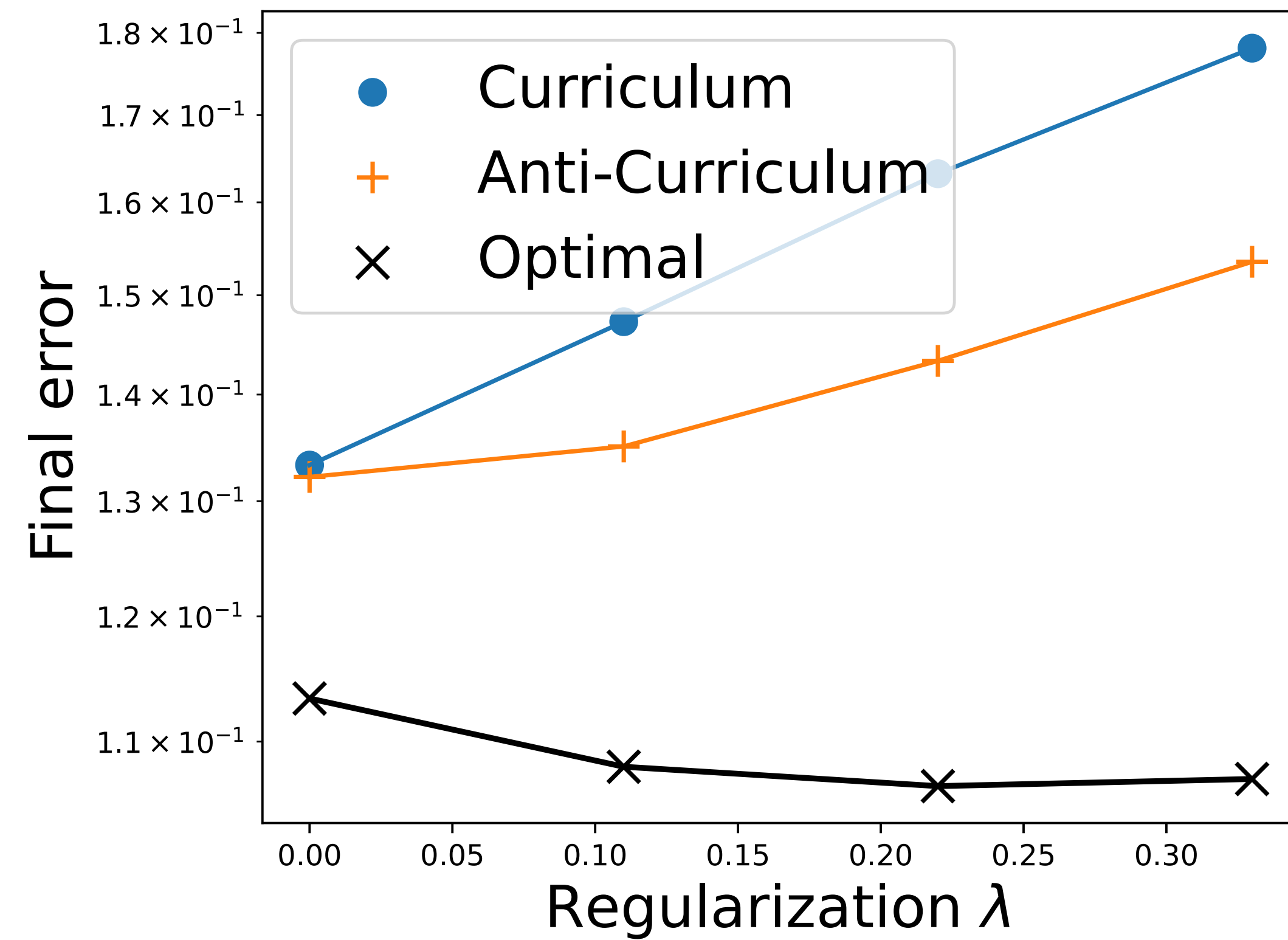
Control: $\mathbf{u} = \Delta$



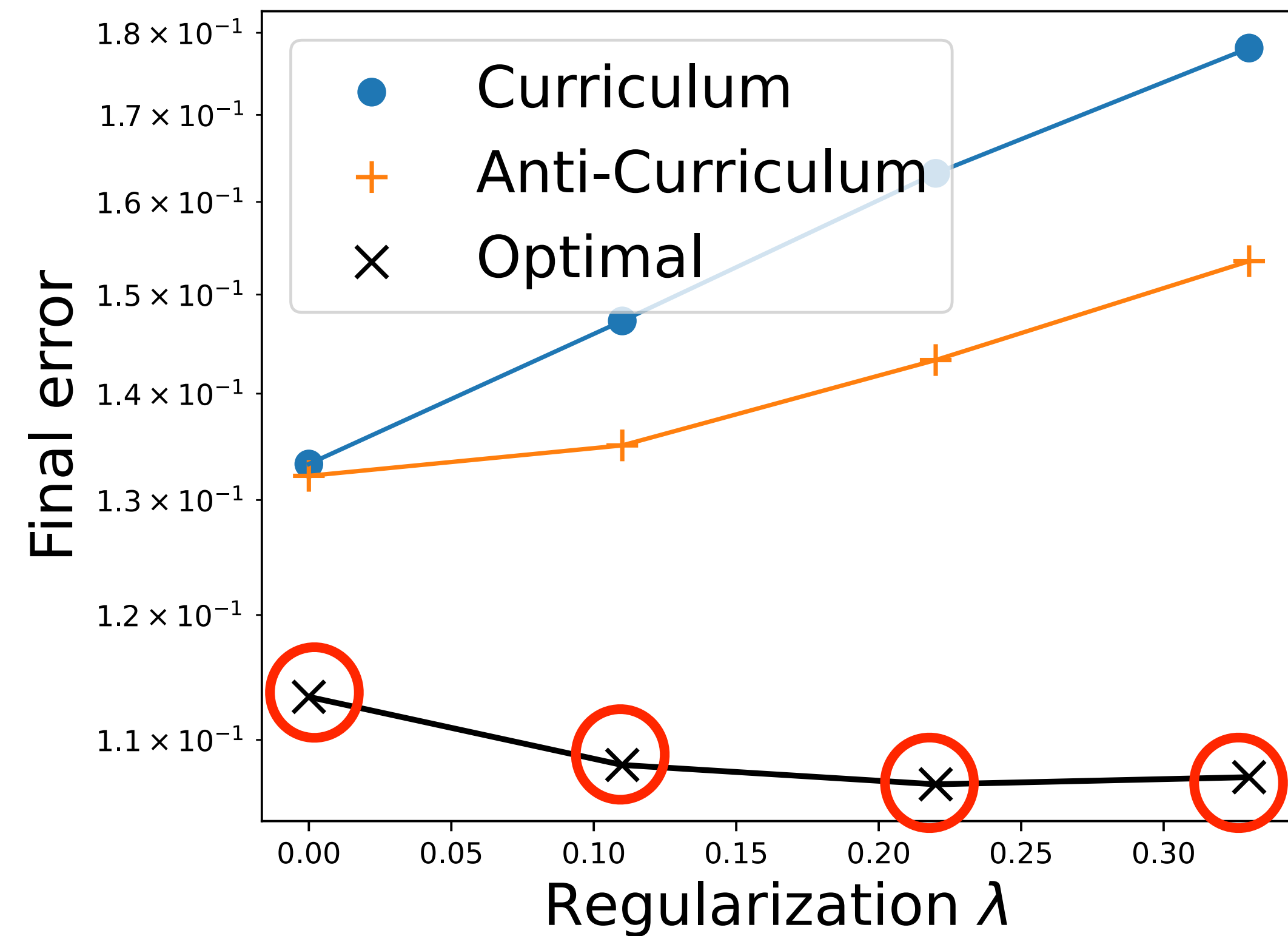
Anti-curriculum
is optimal:



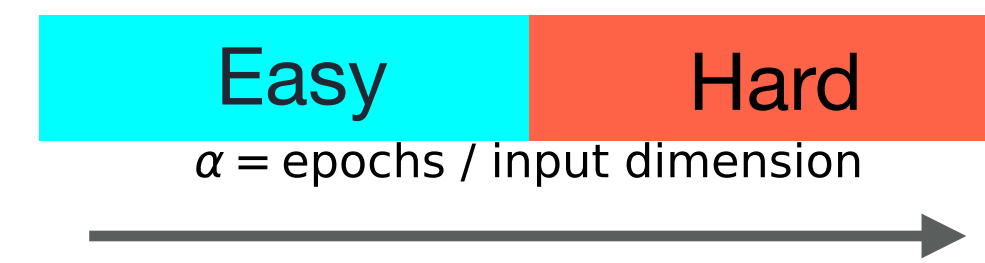
Optimal curriculum protocol



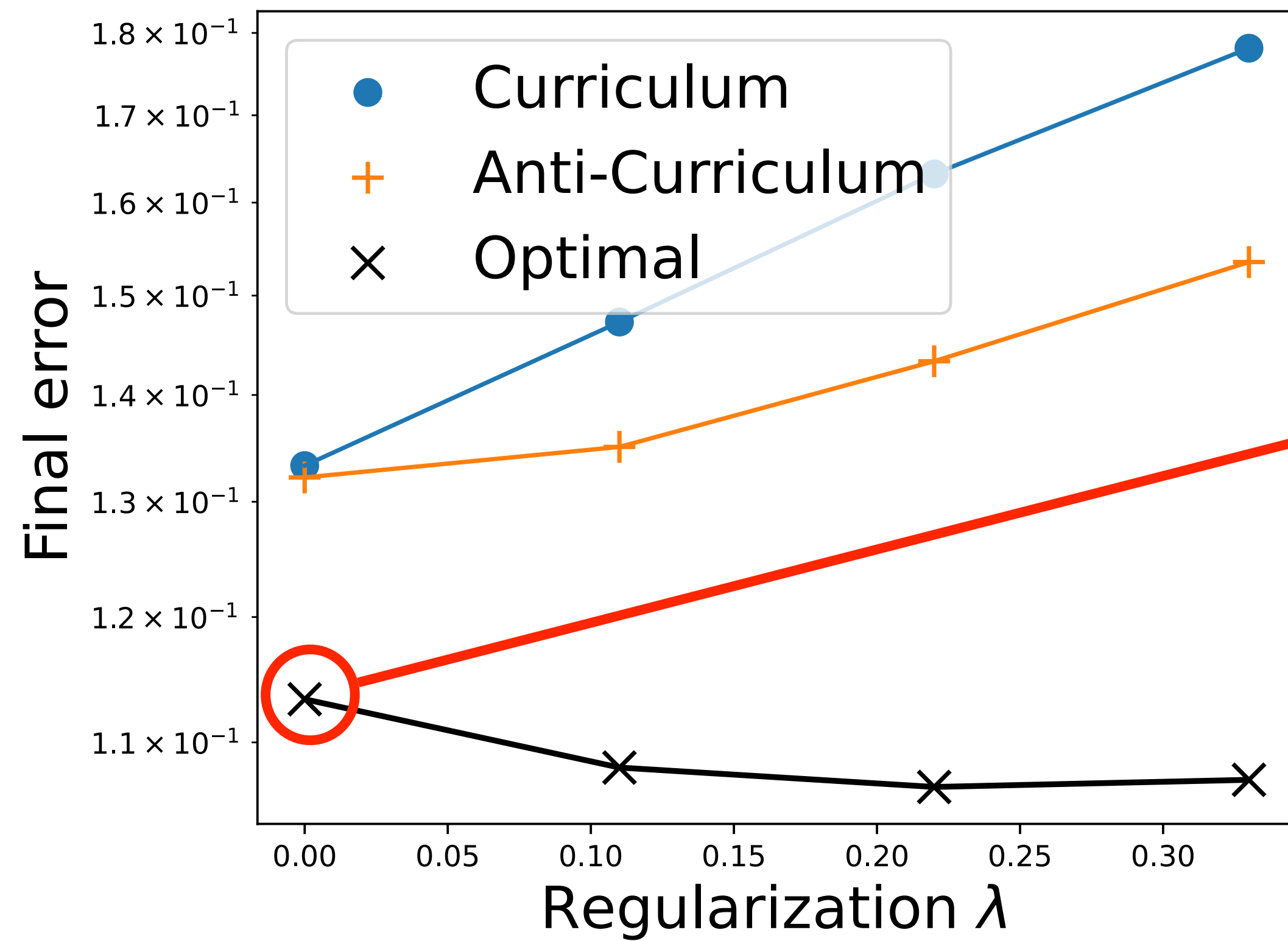
Optimal curriculum protocol



Easy-to-hard
curriculum is optimal:



Optimal curriculum protocol



Easy-to-hard
curriculum is optimal:

