

# The Computational Advantage of Depth: Learning High-Dimensional Hierarchical Functions with Gradient Descent



The Computational Advantage of Depth: Learning High-Dimensional  
Hierarchical Functions with Gradient Descent

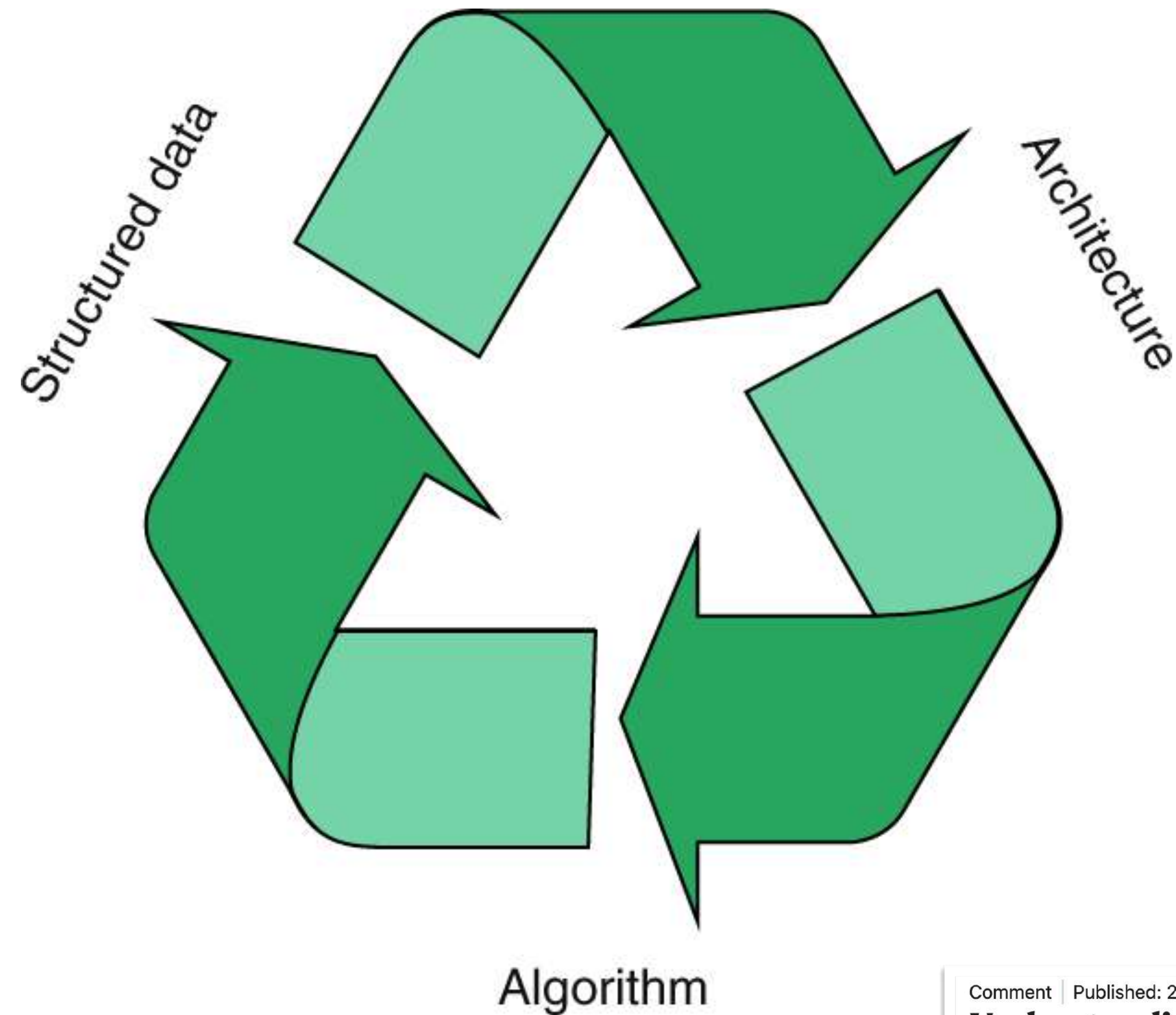
Yatin Dandi<sup>1,2</sup>, Luca Pesce<sup>1</sup>, Lenka Zdeborová<sup>2</sup>, and Florent Krzakala<sup>1</sup>

<sup>1</sup>Ecole Polytechnique Fédérale de Lausanne, Information, Learning and Physics Laboratory. CH-1015 Lausanne, Switzerland.

<sup>2</sup>Ecole Polytechnique Fédérale de Lausanne, Statistical Physics of Computation Laboratory. CH-1015 Lausanne, Switzerland.



# Understanding deep learning is also a job\*



Comment | Published: 26 May 2020

**Understanding deep learning is also a job for physicists**

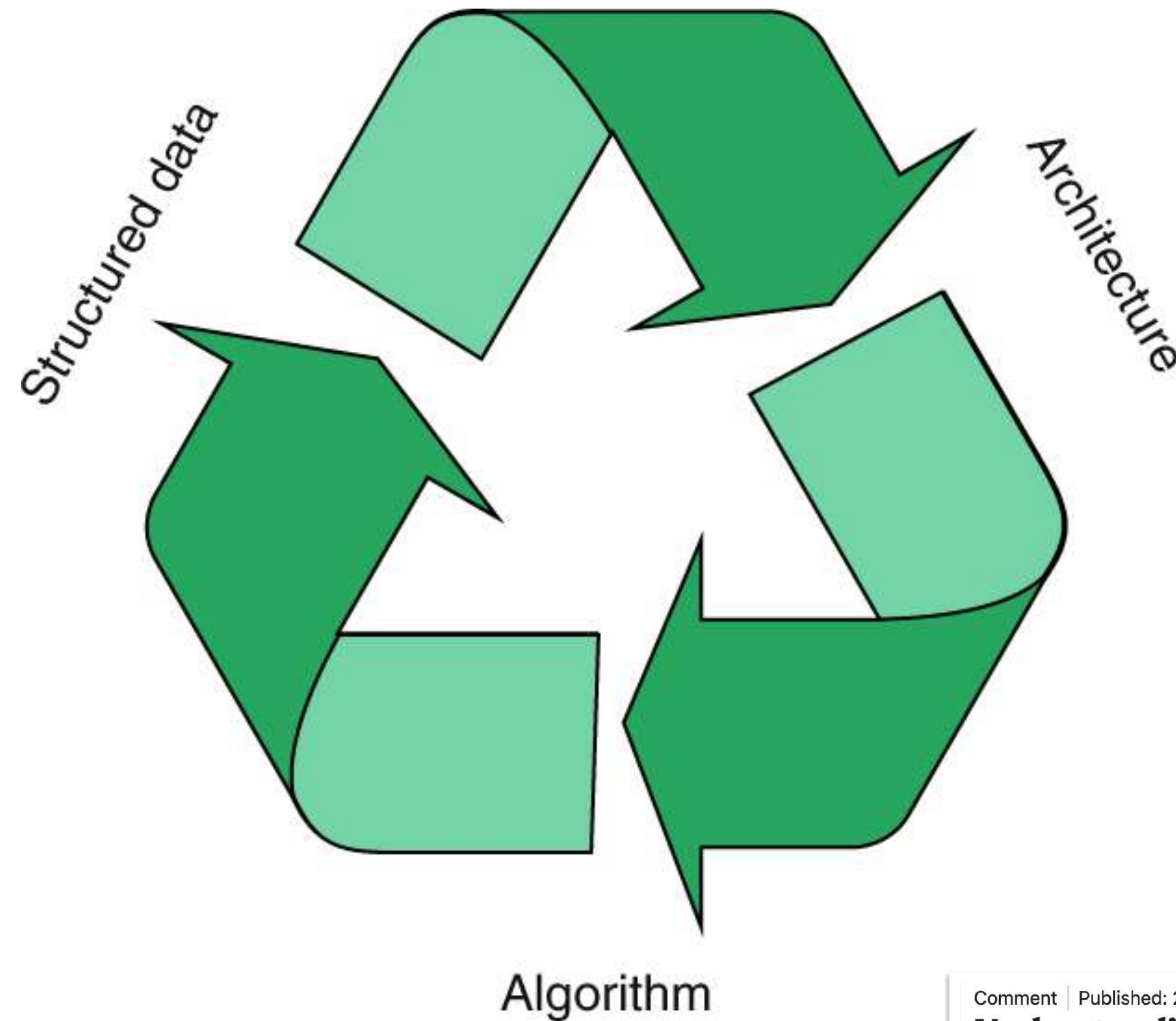
[Lenka Zdeborová](#) 

[Nature Physics](#) **16**, 602–604 (2020) | [Cite this article](#)



# Understanding deep learning is also a job\*

**\*Till AGI**



Comment | Published: 26 May 2020

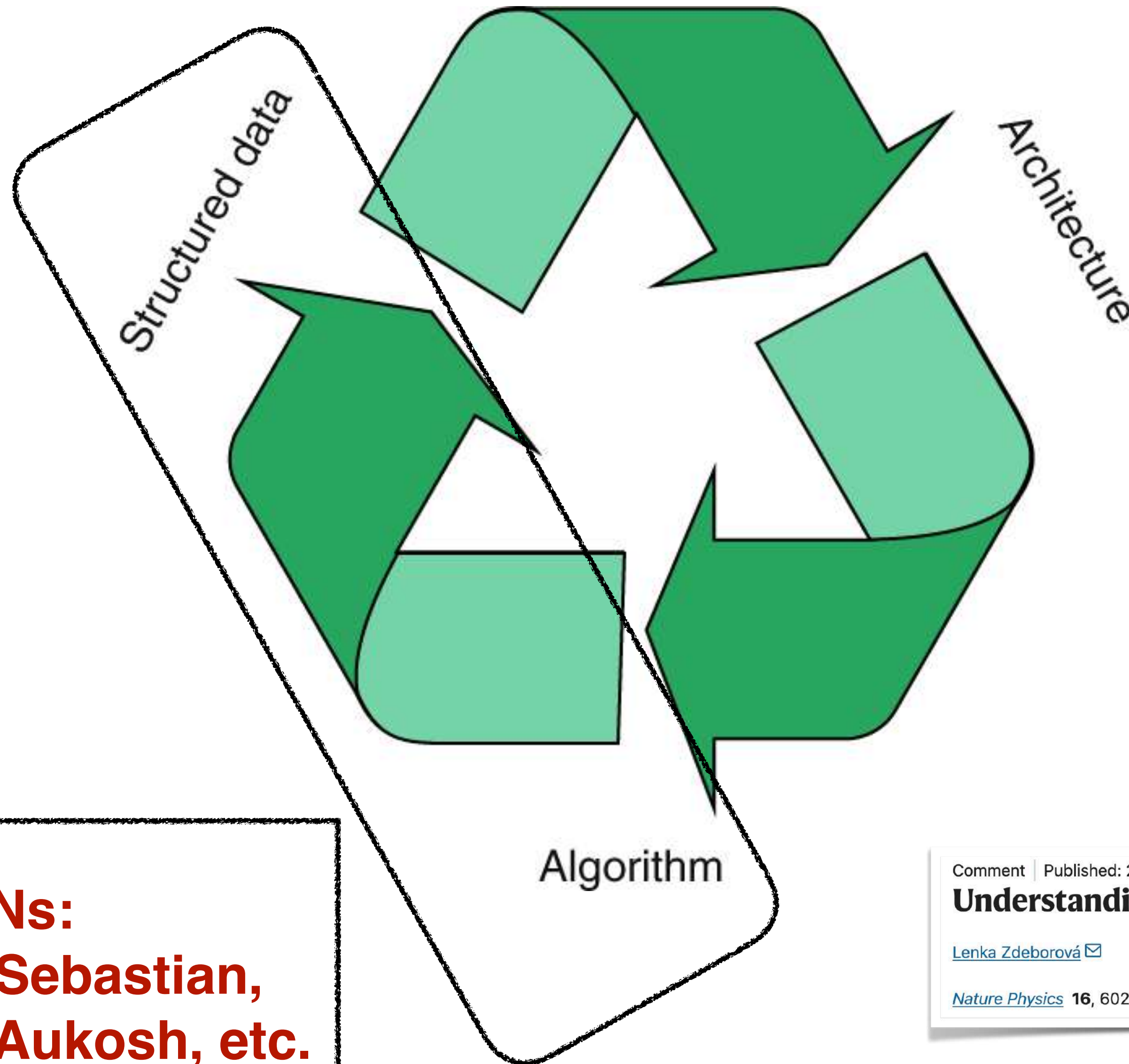
**Understanding deep learning is also a job for physicists**

[Lenka Zdeborová](#) 

[Nature Physics](#) **16**, 602–604 (2020) | [Cite this article](#)

# Understanding deep learning is also a job\*

\*Till AGI



**2-layer NNs:**  
**Bruno's lecture, Sebastian,**  
**Joan, Inbar, Alex, Aukosh, etc.**

Comment | Published: 26 May 2020

**Understanding deep learning is also a job for physicists**

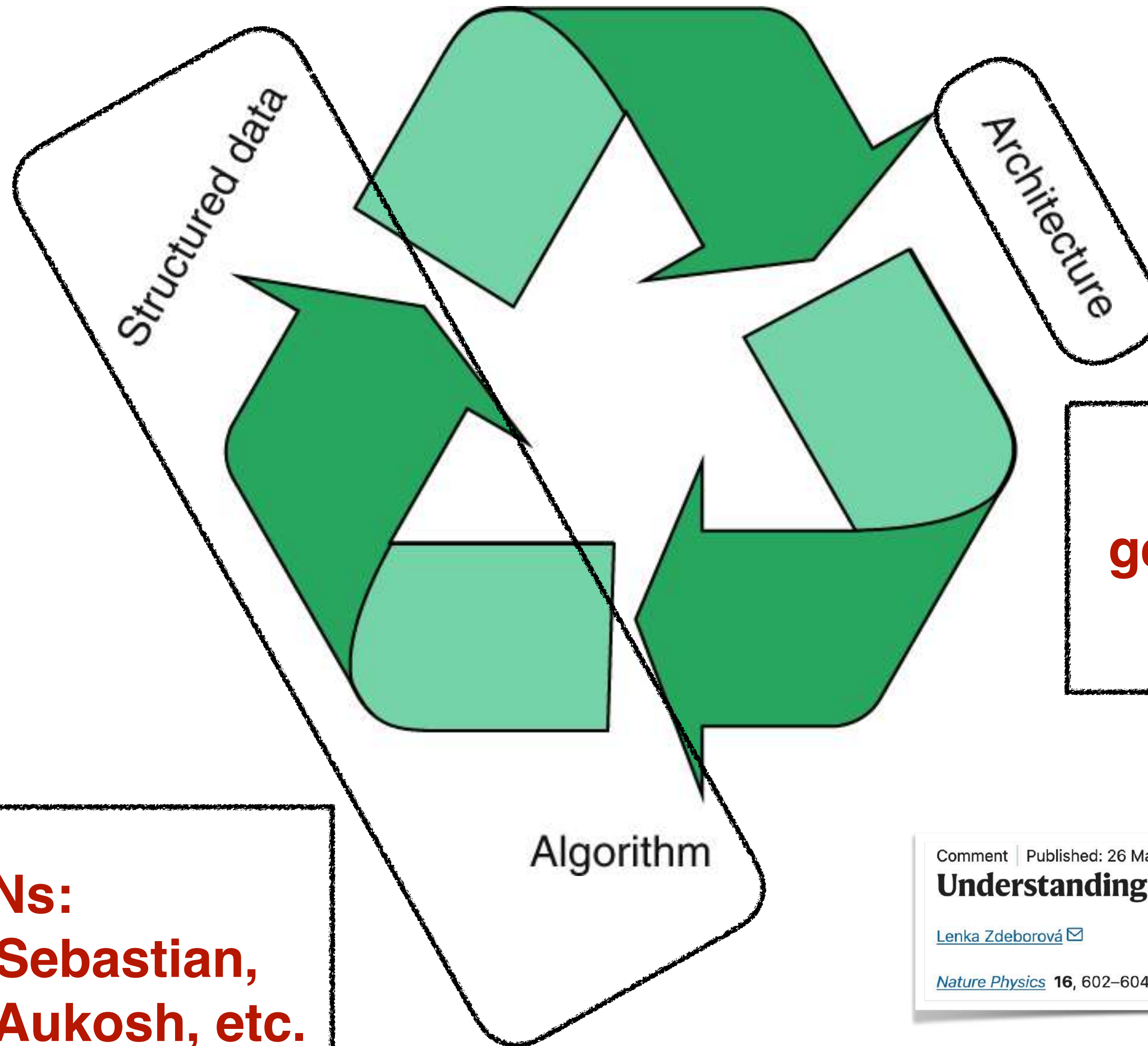
[Lenka Zdeborová](#) 

[Nature Physics](#) **16**, 602–604 (2020) | [Cite this article](#)



# Understanding deep learning is also a job\*

\*Till AGI



**Transformers, CNNs,  
geometric deep learning,  
etc.**

**2-layer NNs:  
Bruno's lecture, Sebastian,  
Joan, Inbar, Alex, Aukosh, etc.**

Comment | Published: 26 May 2020

**Understanding deep learning is also a job for physicists**

[Lenka Zdeborová](#) ✉

[Nature Physics](#) **16**, 602–604 (2020) | [Cite this article](#)

It's called deep learning for a reason.



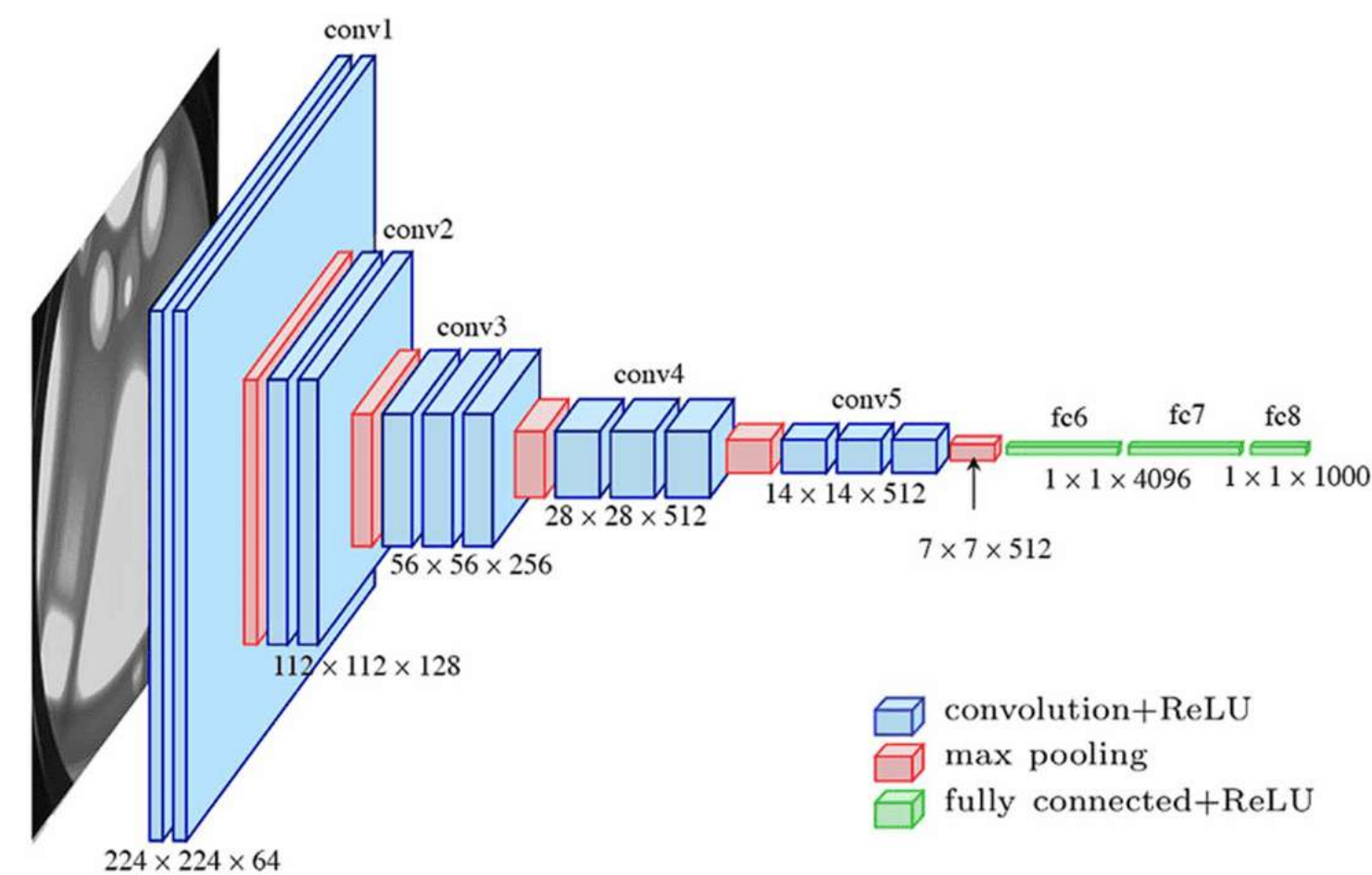
# It's called deep learning for a reason.

VGG

## VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

Karen Simonyan\* & Andrew Zisserman<sup>+</sup>

Visual Geometry Group, Department of Engineering Science, University of Oxford  
{karen, az}@robots.ox.ac.uk



Resnet

by the number of stacked layers (depth). Recent evidence [41, 44] reveals that **network depth is of crucial importance**, and the leading results [41, 44, 13, 16] on the challenging ImageNet dataset [36] all exploit “**very deep**” [41] models, with a depth of sixteen [41] to thirty [16]. Many other non-



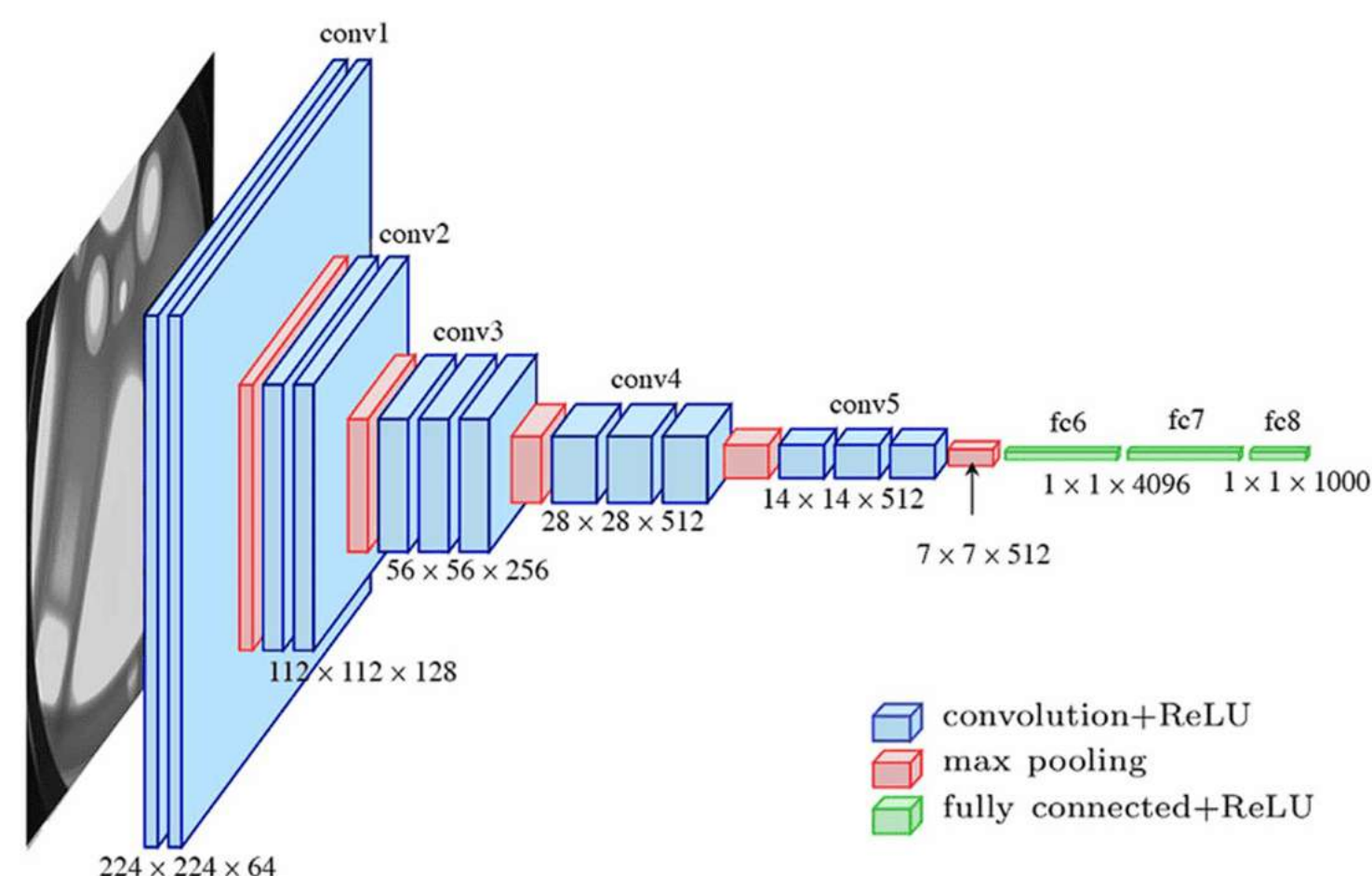
# It's called deep learning for a reason.

VGG

## VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

Karen Simonyan\* & Andrew Zisserman<sup>+</sup>

Visual Geometry Group, Department of Engineering Science, University of Oxford  
{karen,az}@robots.ox.ac.uk



Exponential gains in  
approximation capacity

## The Power of Depth for Feedforward Neural Networks

Ronen Eldan

Weizmann Institute of Science

ronen.eldan@weizmann.ac.il

Ohad Shamir

Weizmann Institute of Science

ohad.shamir@weizmann.ac.il

### Abstract

We show that there is a simple (approximately radial) function on  $\mathbb{R}^d$ , expressible by a small 3-layer feedforward neural networks, which cannot be approximated by any 2-layer network, to more than a certain constant accuracy, unless its width is exponential in the dimension. The result holds for virtually

Resnet

by the number of stacked layers (depth). Recent evidence [41, 44] reveals that **network depth is of crucial importance**, and the leading results [41, 44, 13, 16] on the challenging ImageNet dataset [36] all exploit “**very deep**” [41] models, with a depth of sixteen [41] to thirty [16]. Many other non-

## Deep vs. Shallow Networks: an Approximation Theory Perspective

by

Hrushikesh N. Mhaskar<sup>1</sup> and Tomaso Poggio<sup>2</sup>

1. Department of Mathematics, California Institute of Technology, Pasadena, CA 91125  
Institute of Mathematical Sciences, Claremont Graduate University, Claremont, CA 91711.  
hrushikesh.mhaskar@cgu.edu

2. Center for Brains, Minds, and Machines, McGovern Institute for Brain Research,  
Massachusetts Institute of Technology, Cambridge, MA, 02139.  
tp@mit.edu

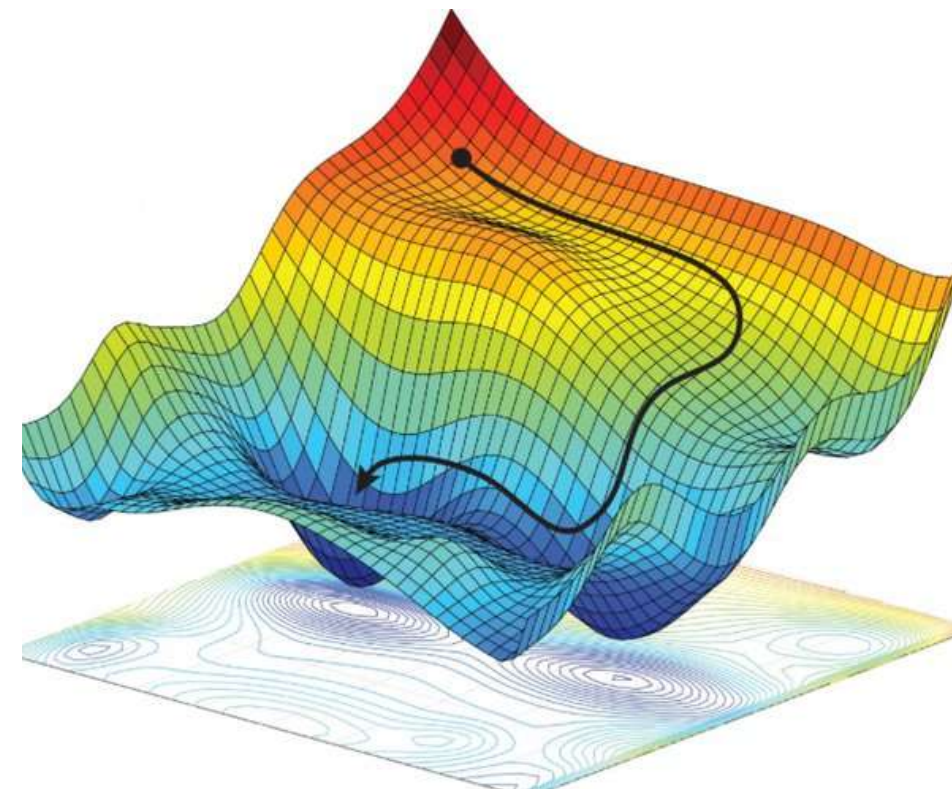


# What do we expect to show?

# What do we expect to show?

## The quest for adaptivity

Posted on June 17, 2021 by Francis Bach



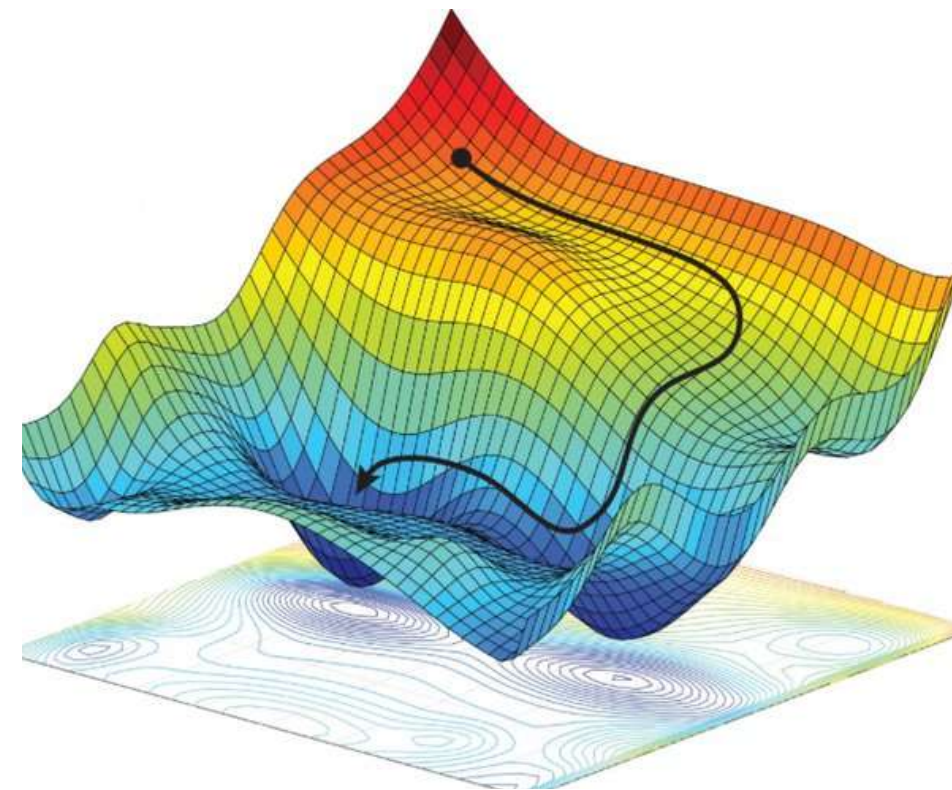


# What do we expect to show?

## The quest for adaptivity

Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**

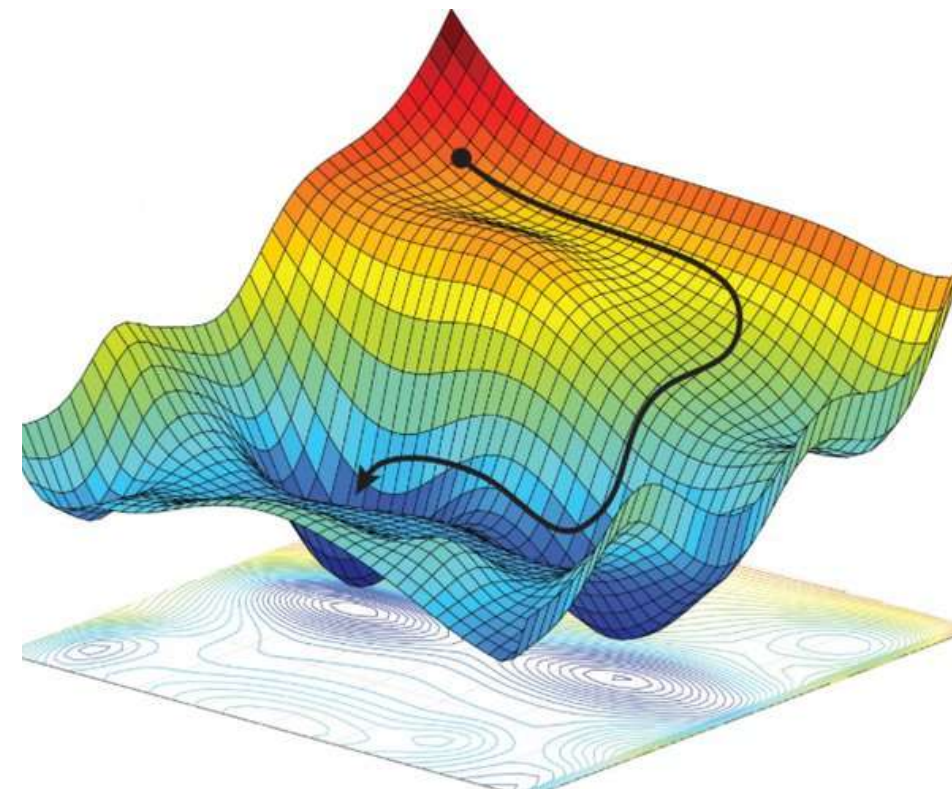
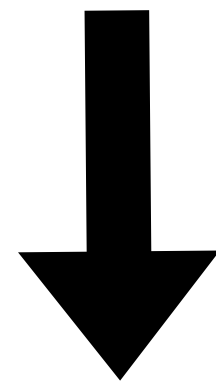


# What do we expect to show?

## The quest for adaptivity

Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**



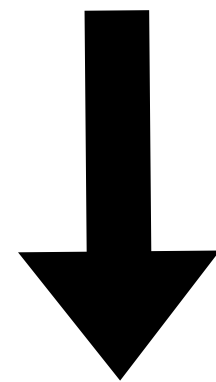


# What do we expect to show?

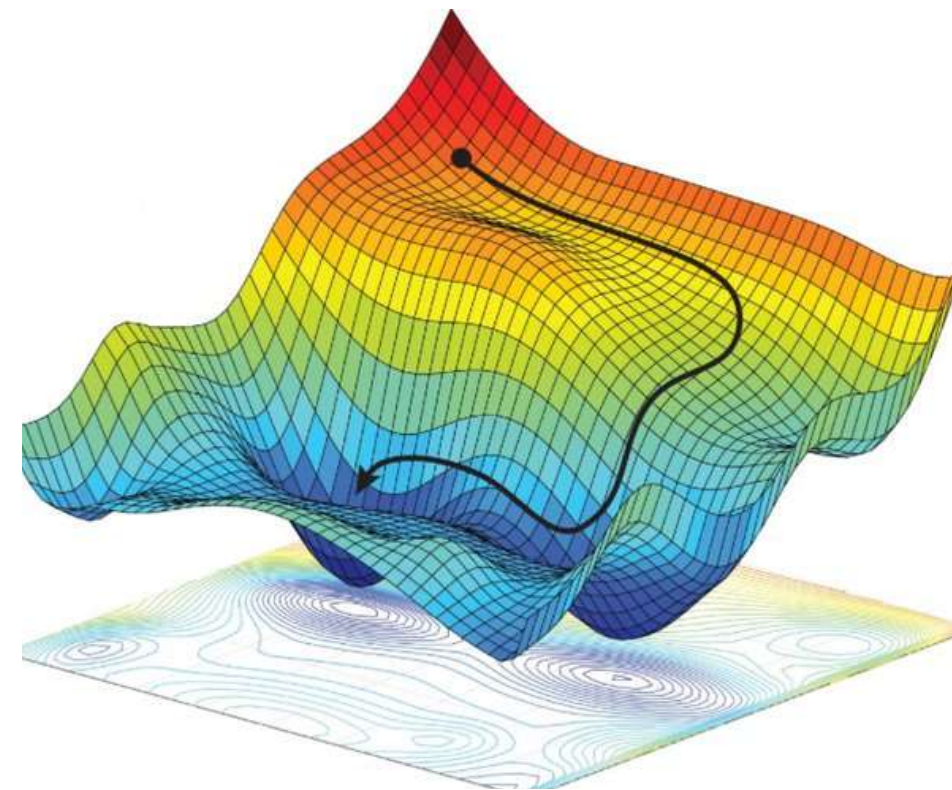
## The quest for adaptivity

Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**



**Lower-sample Complexity  
with adaptive algorithms**

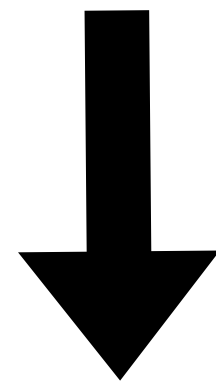


# What do we expect to show?

## The quest for adaptivity

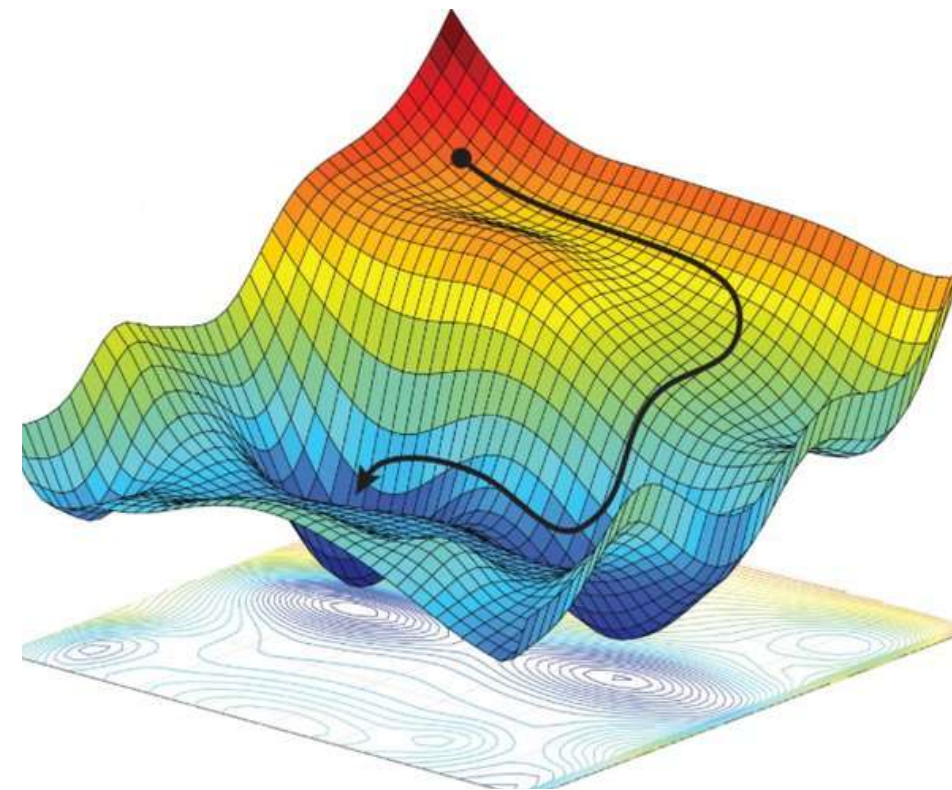
Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**



**Lower-sample Complexity  
with adaptive algorithms**

**Gaussian data+ single, multi-  
index models**



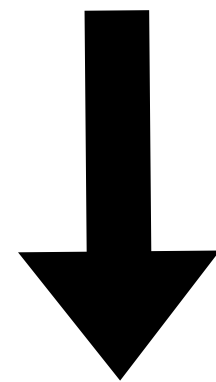


# What do we expect to show?

## The quest for adaptivity

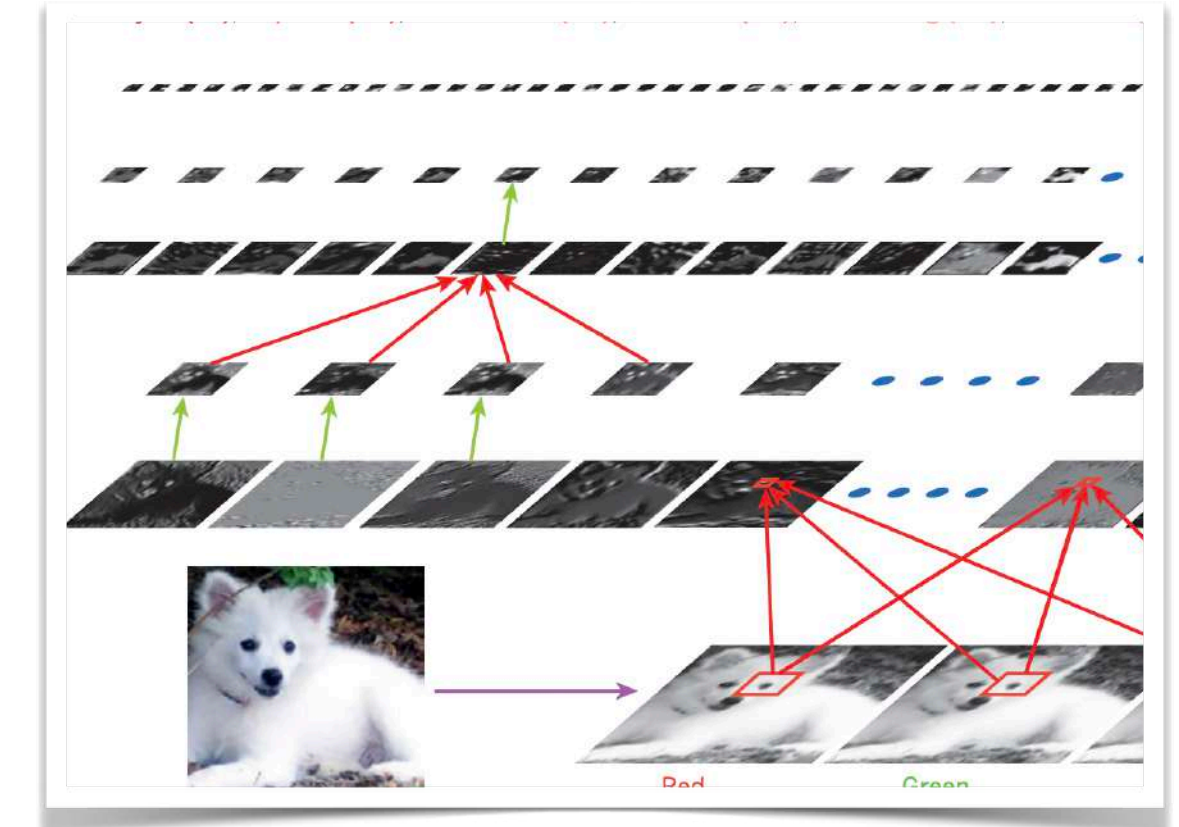
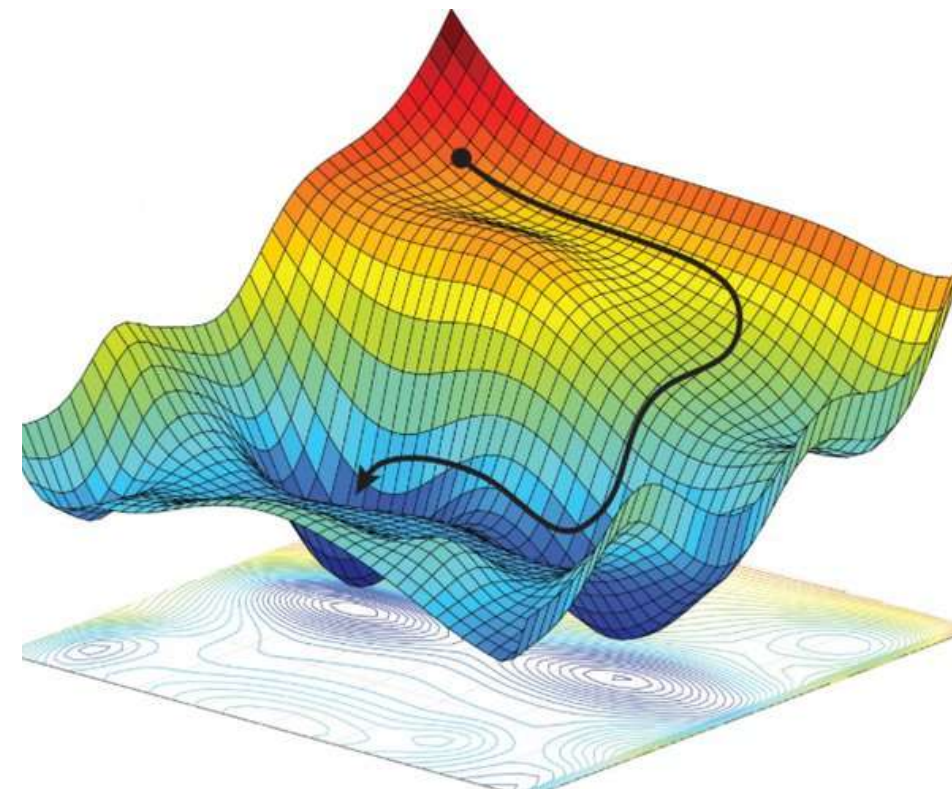
Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**



**Lower-sample Complexity  
with adaptive algorithms**

**Gaussian data+ single, multi-  
index models**



**Hierarchical structure**

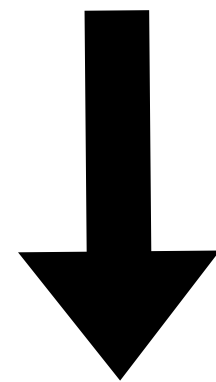


# What do we expect to show?

## The quest for adaptivity

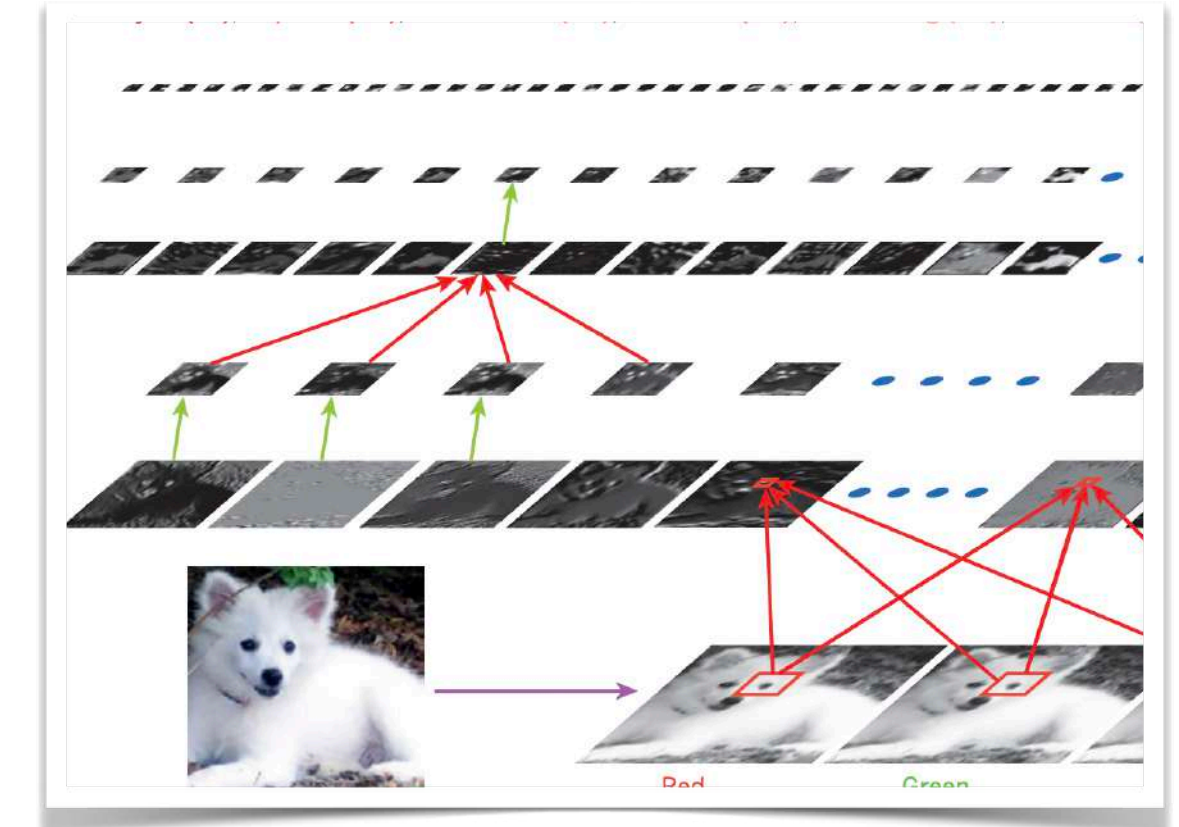
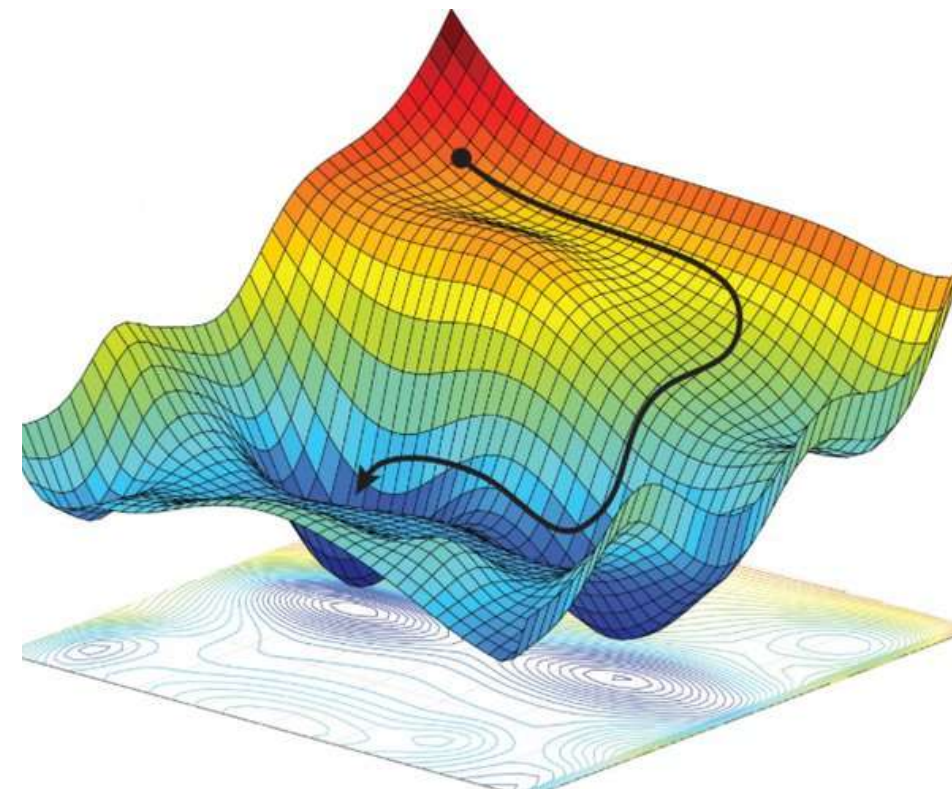
Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**

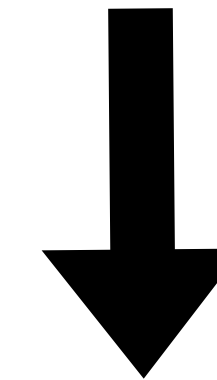


**Lower-sample Complexity  
with adaptive algorithms**

**Gaussian data+ single, multi-  
index models**



**Hierarchical structure**



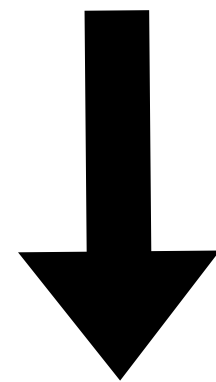


# What do we expect to show?

## The quest for adaptivity

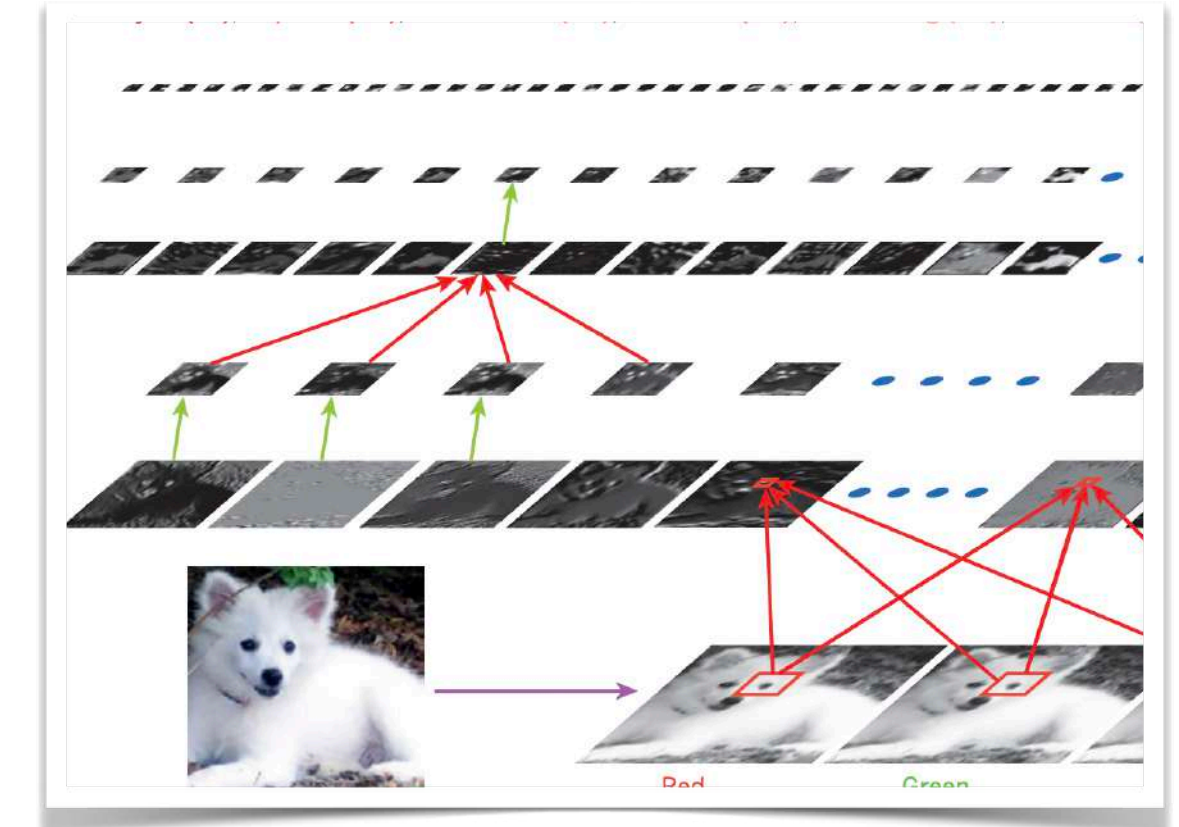
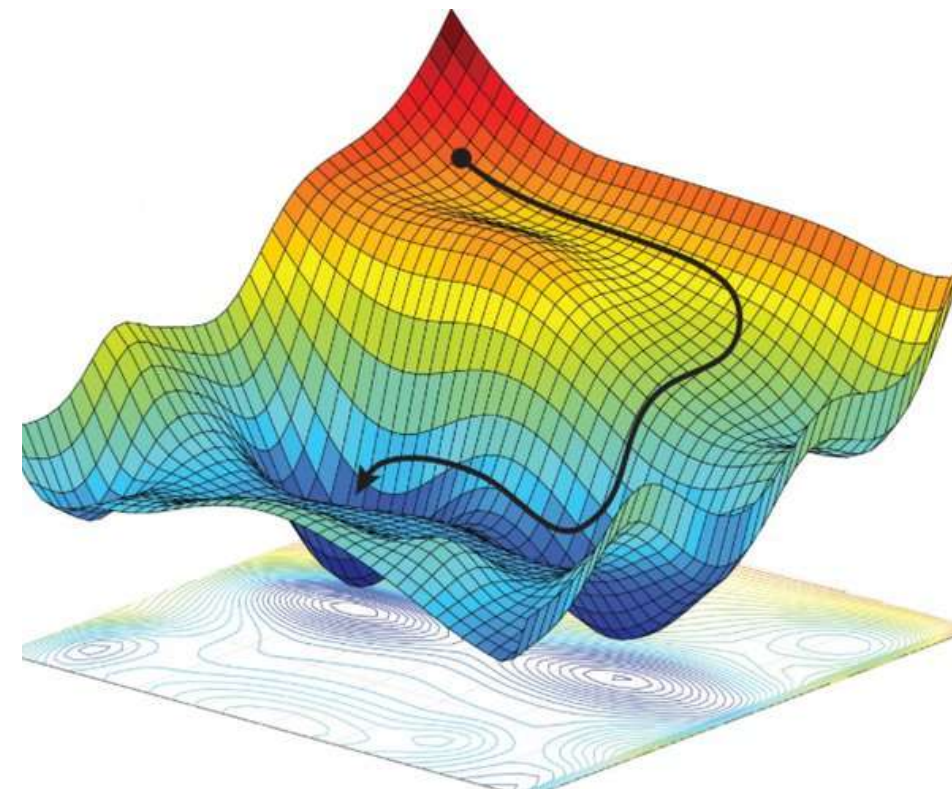
Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**

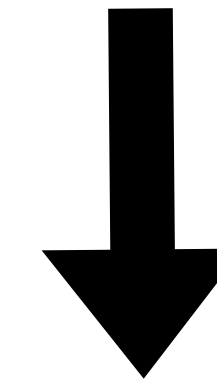


**Lower-sample Complexity  
with adaptive algorithms**

**Gaussian data+ single, multi-  
index models**



**Hierarchical structure**



**Lower-sample Complexity  
With deep/multiple levels  
of adaptation.**

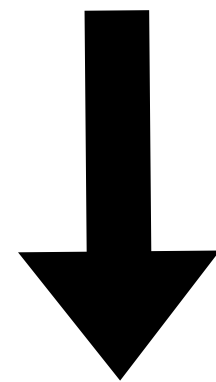


# What do we expect to show?

## The quest for adaptivity

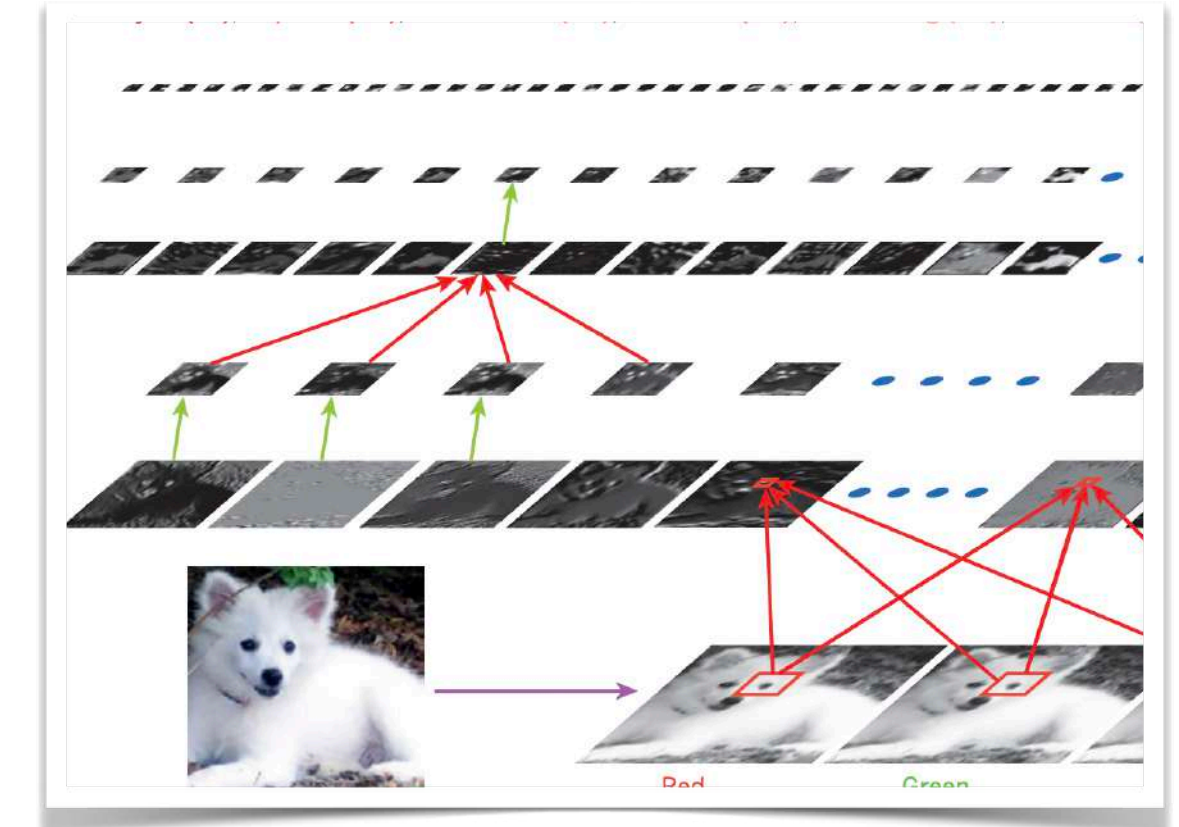
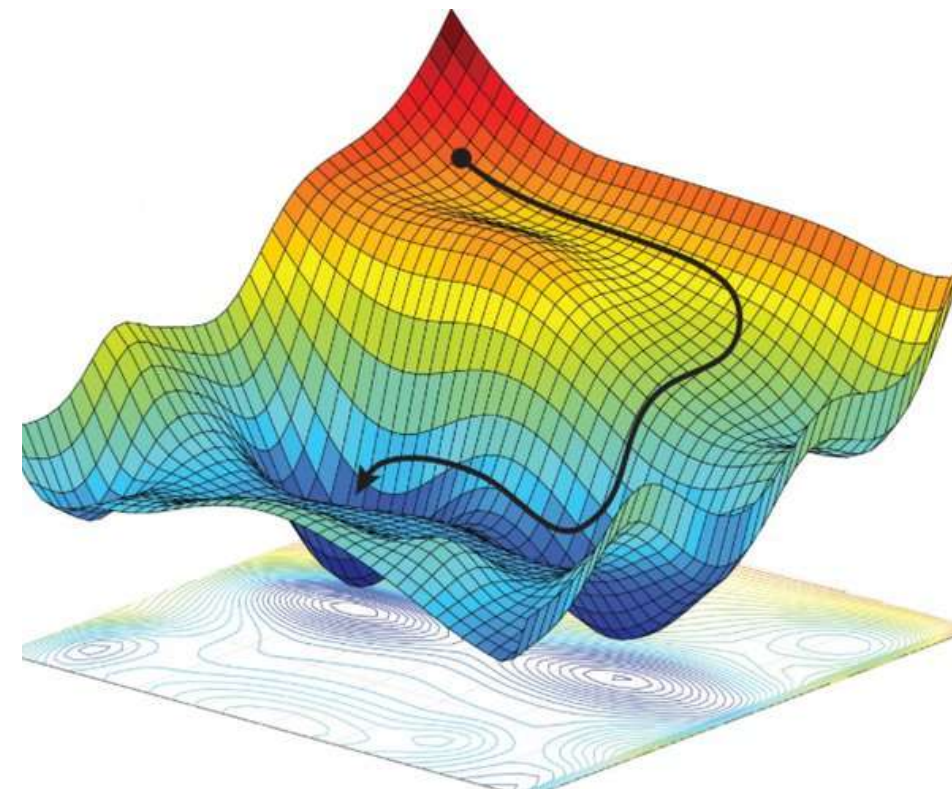
Posted on June 17, 2021 by Francis Bach

**Low-dimensional  
structure**

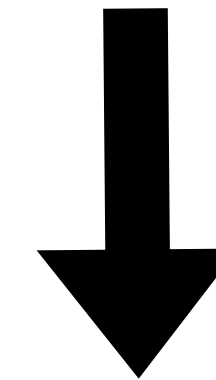


**Lower-sample Complexity  
with adaptive algorithms**

**Gaussian data+ single, multi-  
index models**



**Hierarchical structure**



**Lower-sample Complexity  
With deep/multiple levels  
of adaptation.**

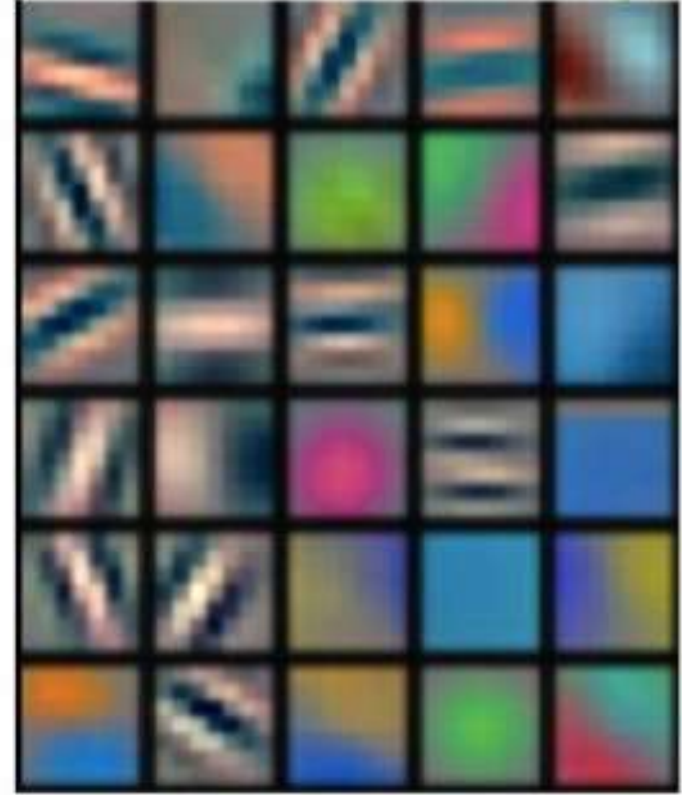
**?**



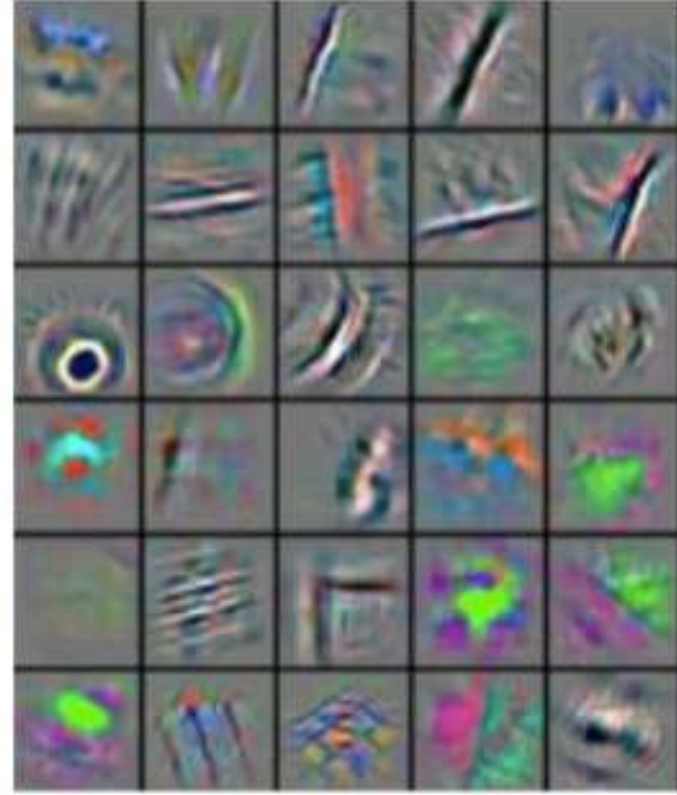
# What's “hierarchical”?

# What's “hierarchical”?

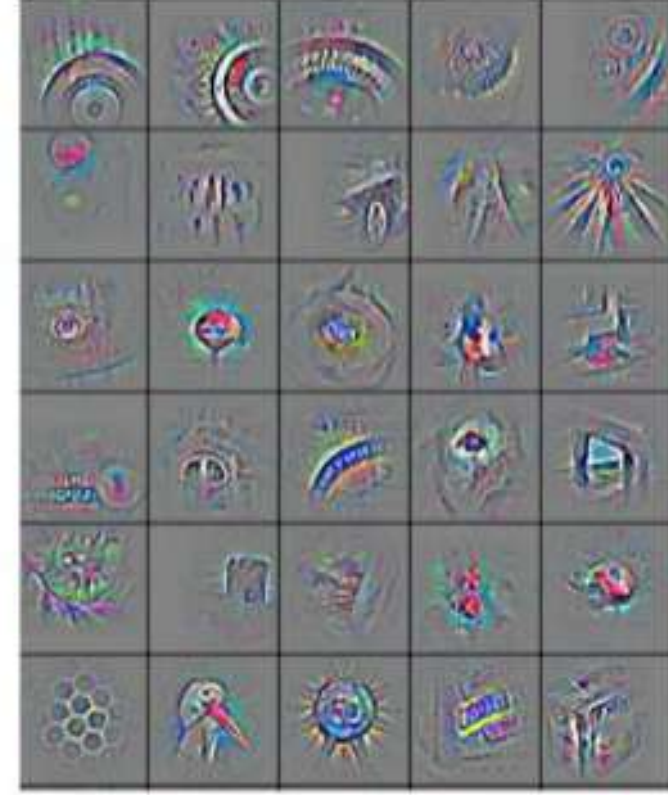
low-level features



mid-level features



high-level features



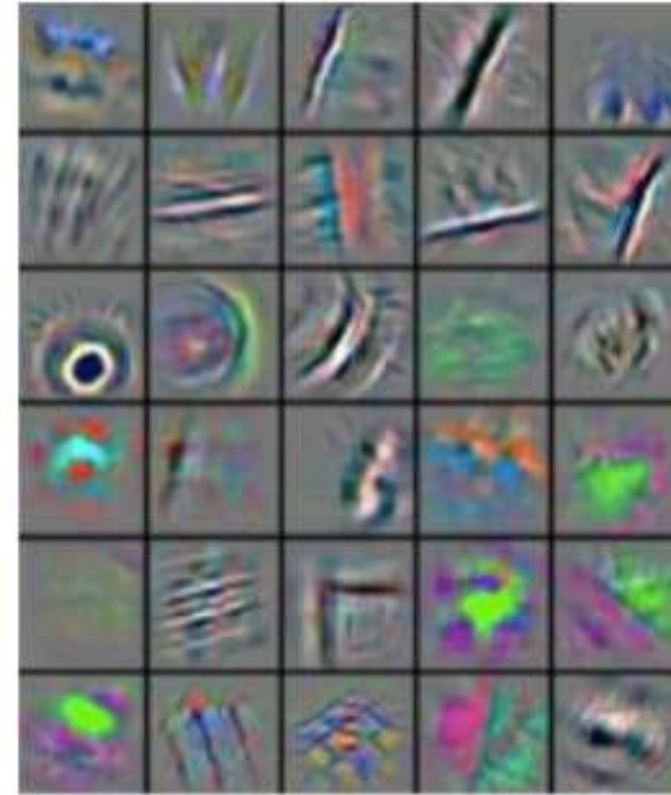
# What's “hierarchical”?

**Large search  
space for target**

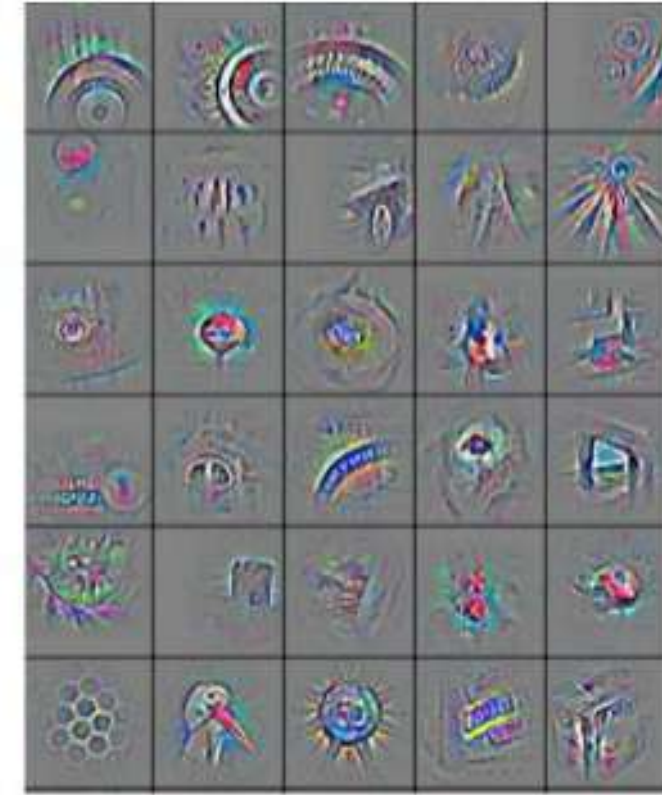
low-level features



mid-level features



high-level features



**Large  
effective dimension**



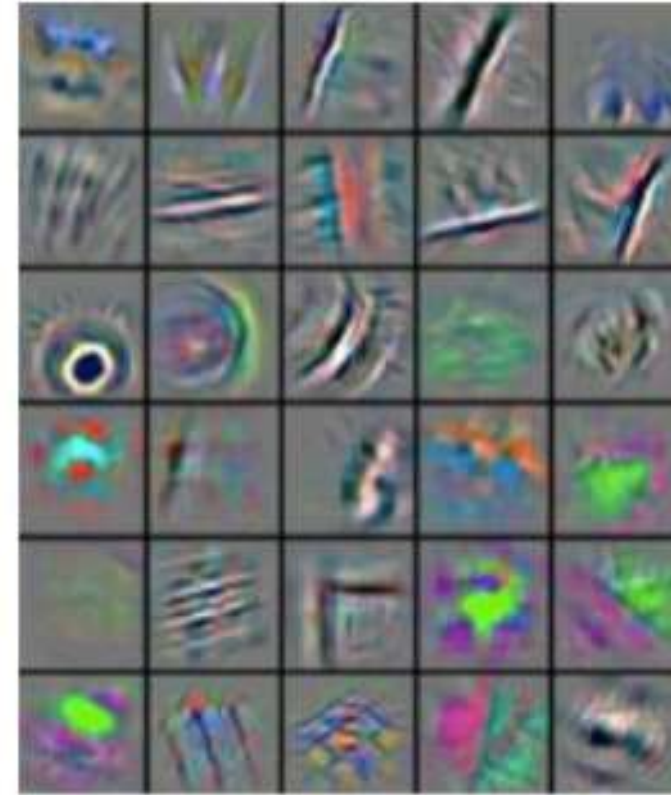
# What's “hierarchical”?

**Large search  
space for target**

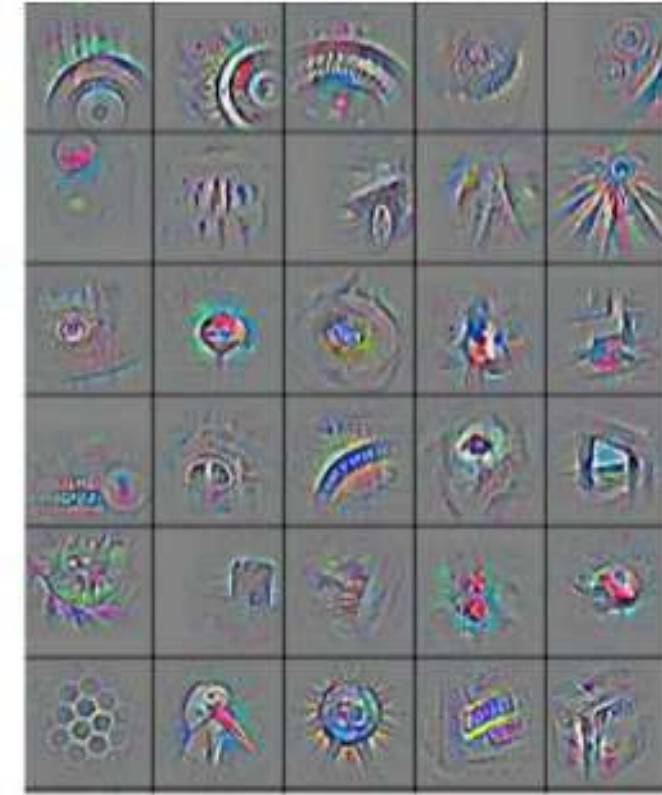
low-level features



mid-level features



high-level features



**Large  
effective dimension**

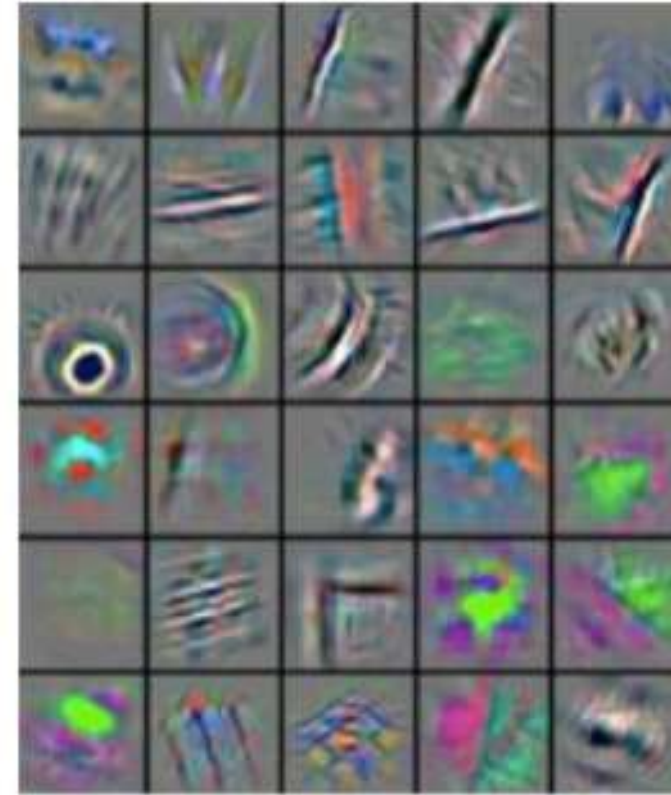


# What's “hierarchical”?

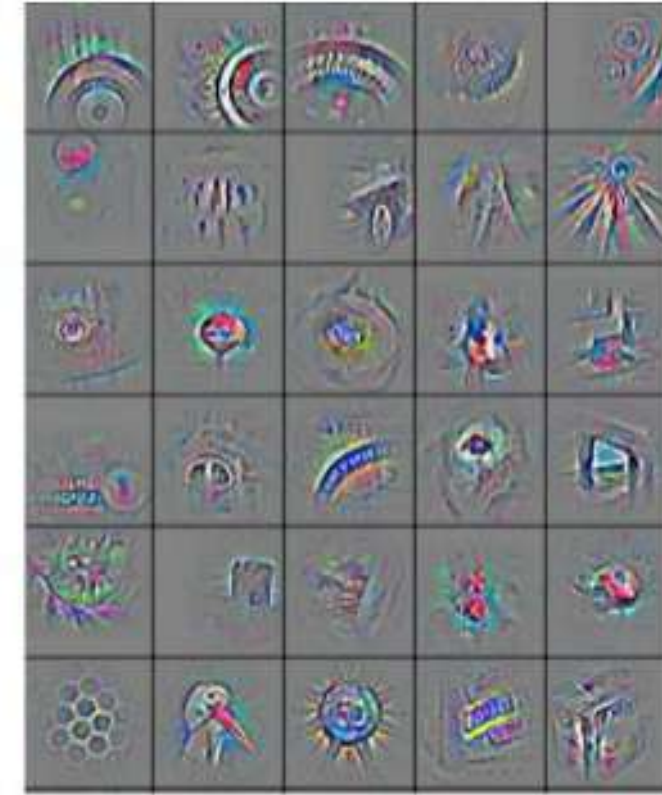
low-level features



mid-level features



high-level features



**Large search  
space for target**

**Small search space  
for target**

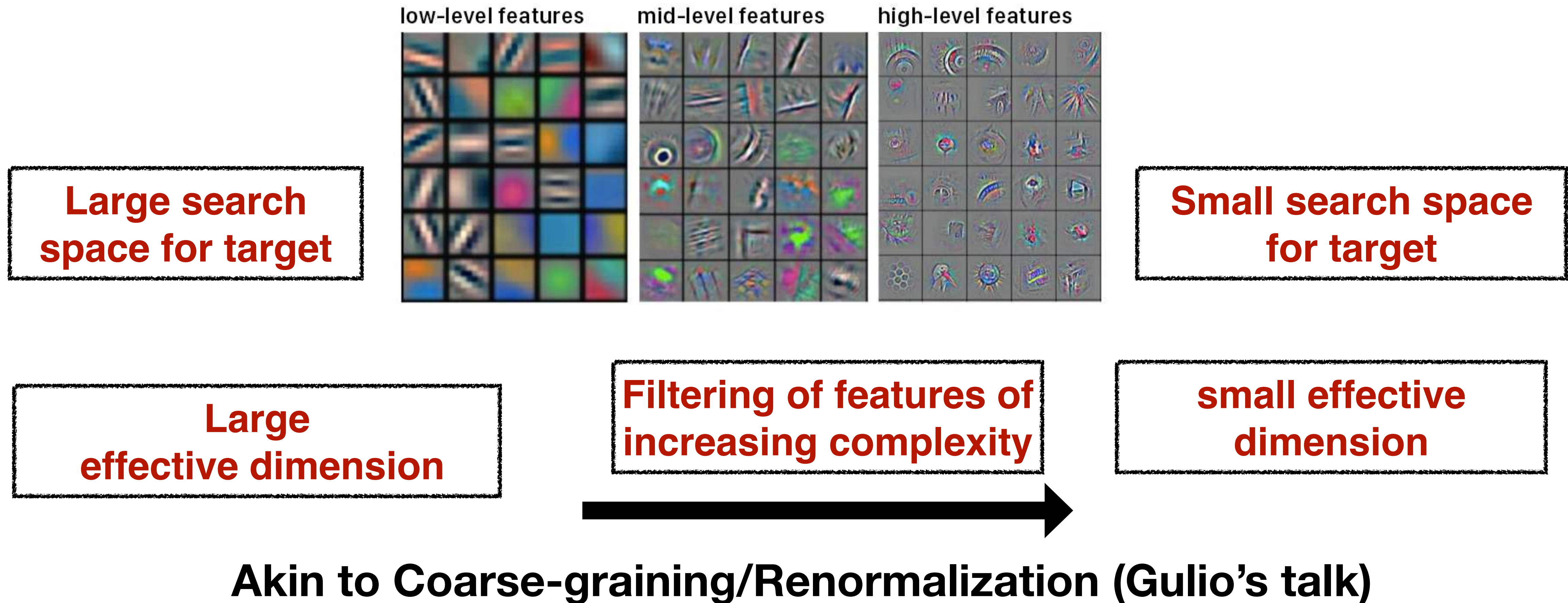
**Large  
effective dimension**

**small effective  
dimension**





# What's “hierarchical”?



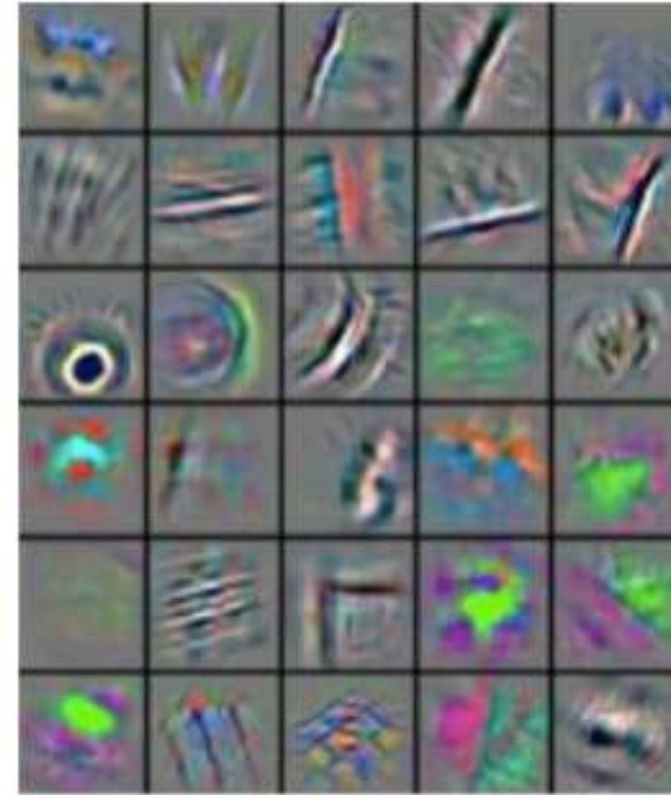


# What's “hierarchical”?

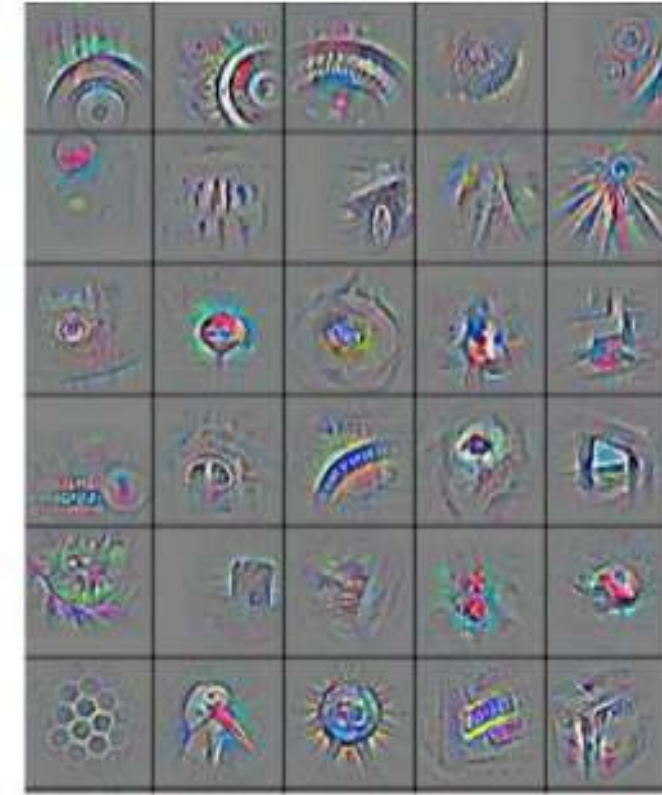
low-level features



mid-level features



high-level features



**Large search  
space for target**

**See also Cagnetta et al.  
2025, Mossel et al., 2016,  
Bruna & Mallat 2013**

**Small search space  
for target**

**Large  
effective dimension**

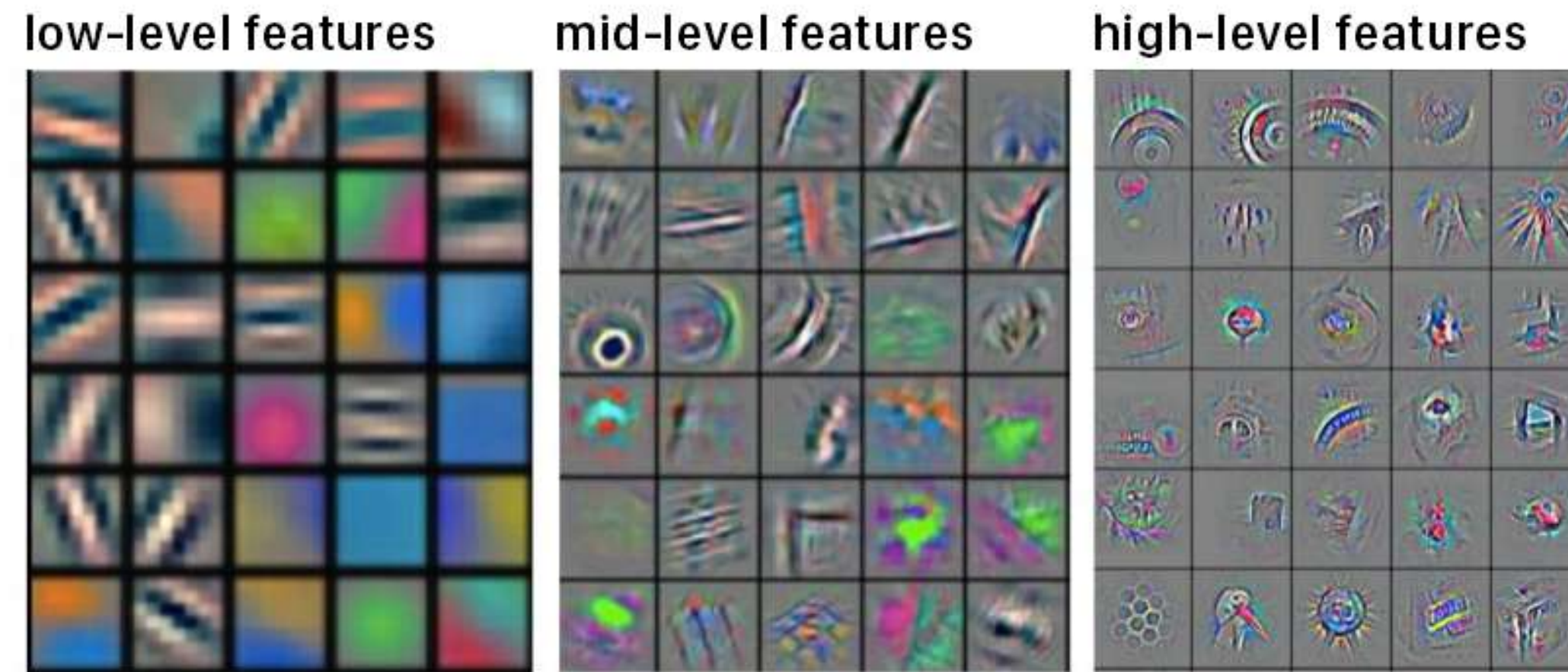
**Filtering of features of  
increasing complexity**

**small effective  
dimension**

**Akin to Coarse-graining/Renormalization (Gulio's talk)**



# What's “hierarchical”?



**Large search  
space for target**

**See also Cagnetta et al.  
2025, Mossel et al., 2016,  
Bruna & Mallat 2013**

**Small search space  
for target**

**Large  
effective dimension**

**Filtering of features of  
increasing complexity**

**small effective  
dimension**

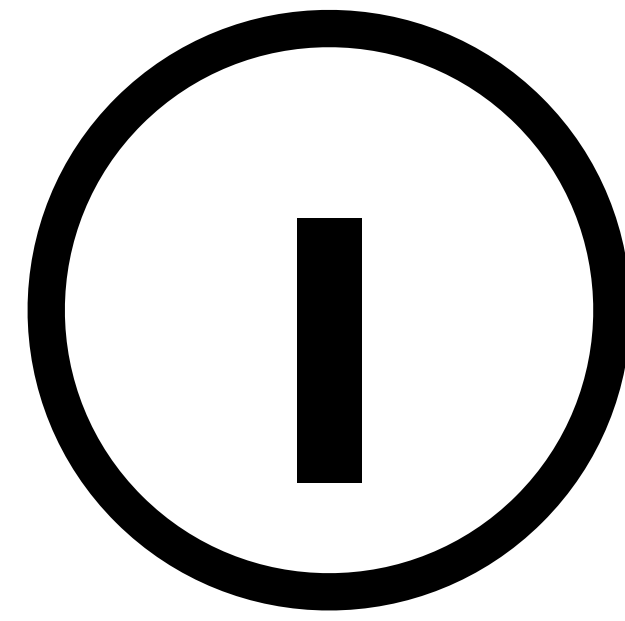


**Akin to Coarse-graining/Renormalization (Gulio's talk)**

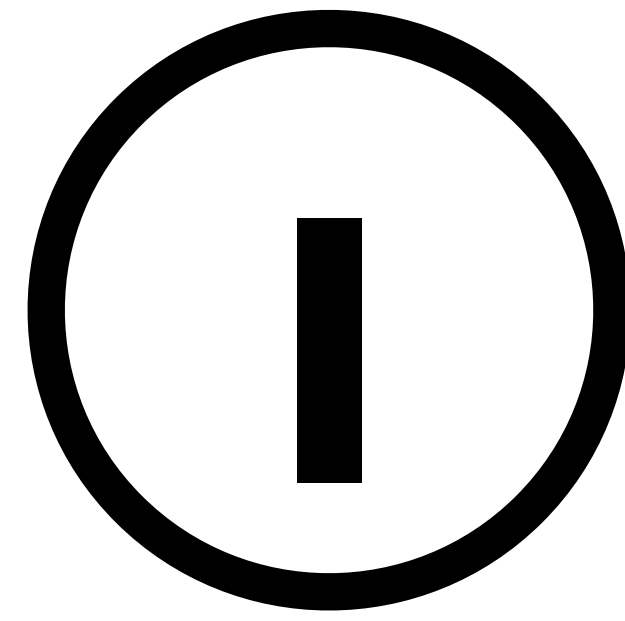
**Can we understand this in some  
analyzable setting?**







# **Recap of two-layer NNs for AGI**



# Recap of two-layer NNs for AGI



**Jason Lee** ✓

@jasondeanlee



At the [@SimonsInstitute](#) working on AGI (Artificial Gaussian Intelligence)

5:12 PM · Nov 15, 2024 · **11.5K** Views



# Neural Networks without Feature Learning

## Random features

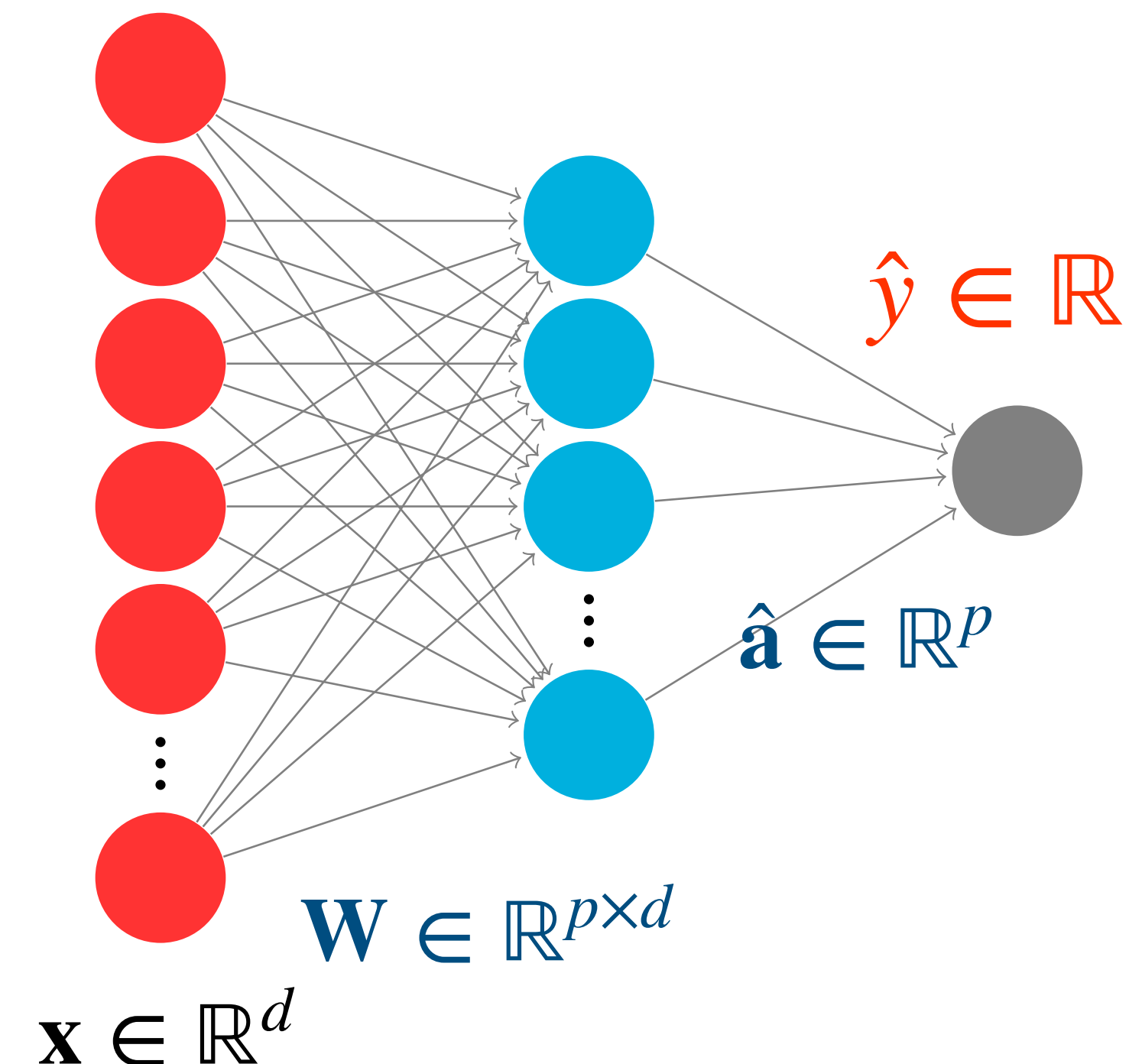
[Balcan, Blum, Vempala '06, Rahimi-Recht '17...]

No training of the first layer:  $W$  is fixed

$$\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=1}^p \hat{a}_i \sigma_i(\langle \mathbf{w}_i, \mathbf{x} \rangle) = \sum_{i=1}^p \hat{a}_i \Phi_{\text{CK}}(\mathbf{x})$$

Computationally easy (linear regression)

$$\hat{y} = \hat{f}(\mathbf{x}) = \hat{\mathbf{a}} \cdot \sigma(W\mathbf{x})$$



Very popular setting among theoreticians

Equivalent to Neural Tangent Kernel/Lazy Regime/Kernel methods/ etc..

[Jacot, Gabriel, Hongler '18; Lee, Jaehoon, et al. 18; Chizat, Bach '19,...]

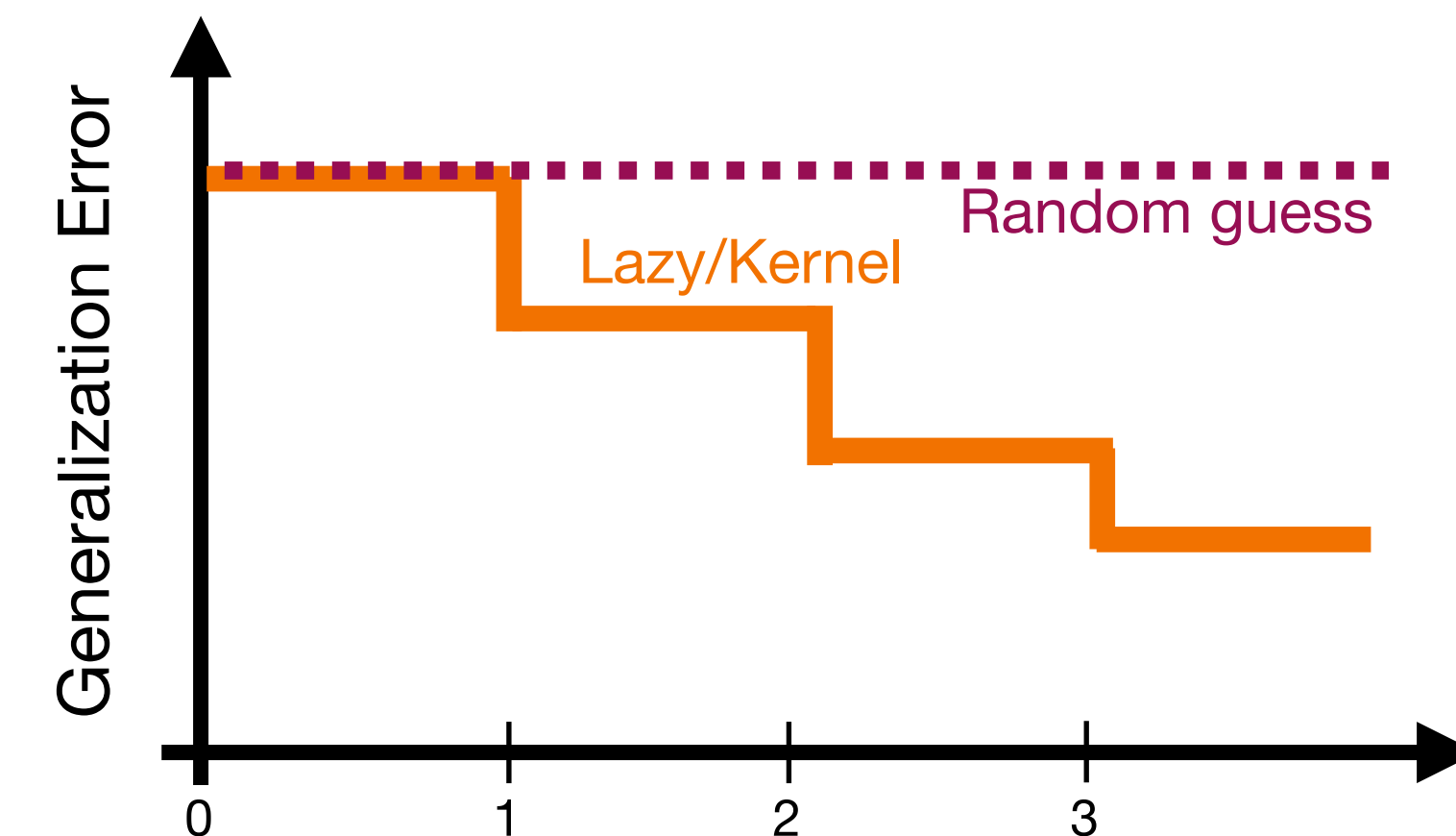
# Recap of Sample Complexity

**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$



See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]



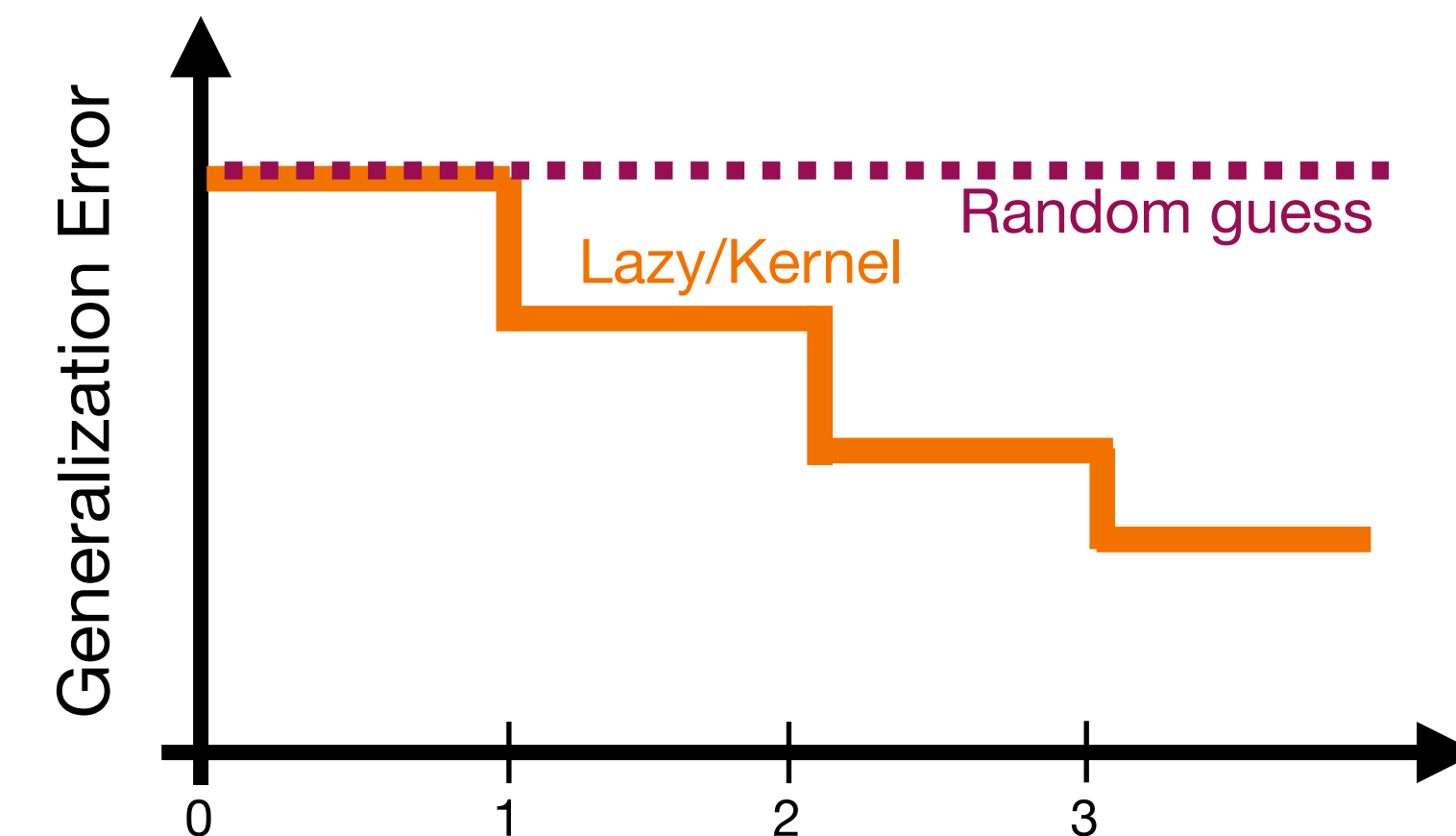
# Recap of Sample Complexity

**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$
$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$

$(n, p) = O(d)$



See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]

# Recap of Sample Complexity

**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

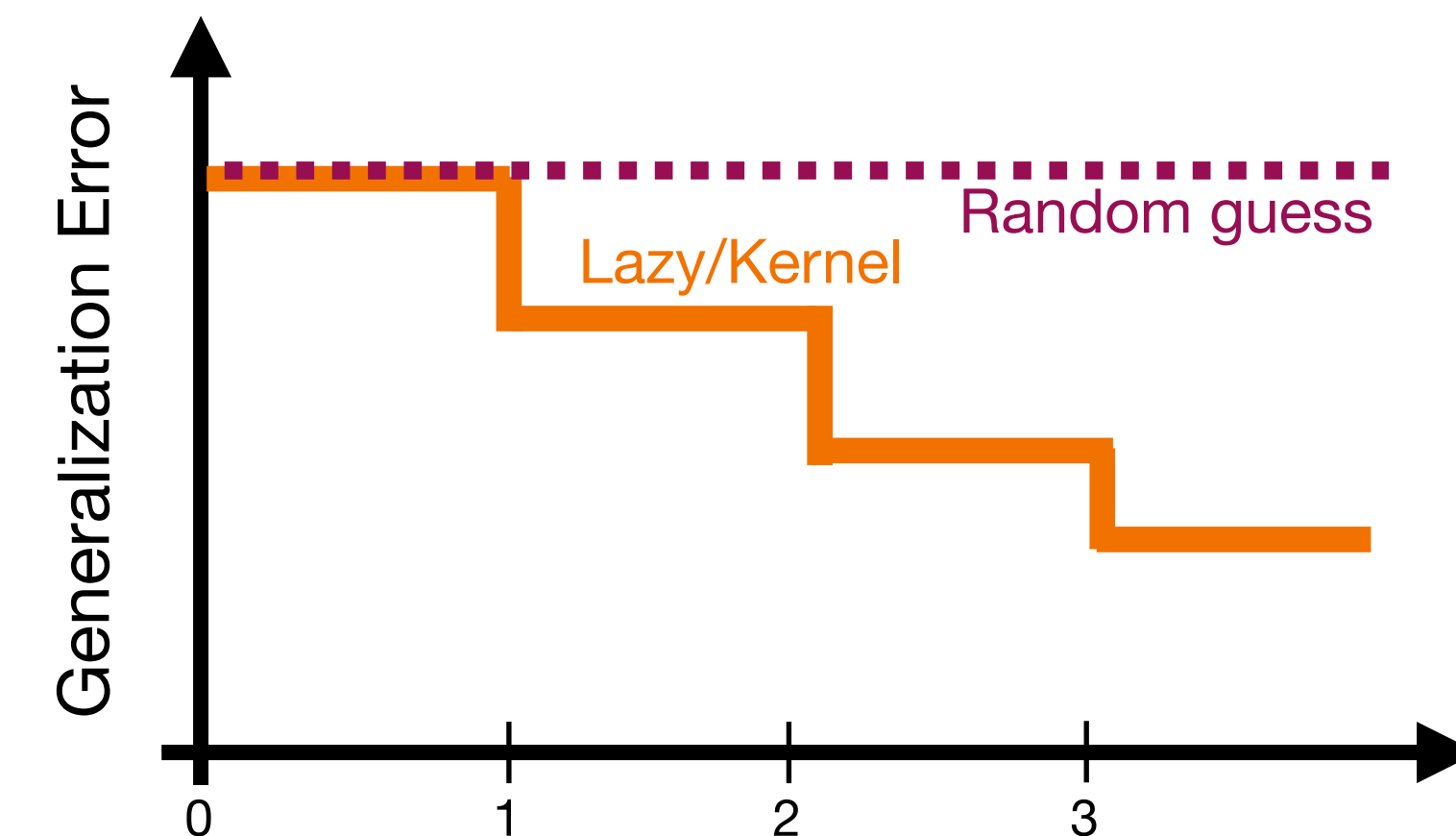
In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$

$$(n, p) = O(d)$$

$$(n, p) = O(d^2)$$



See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]



# Recap of Sample Complexity

**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

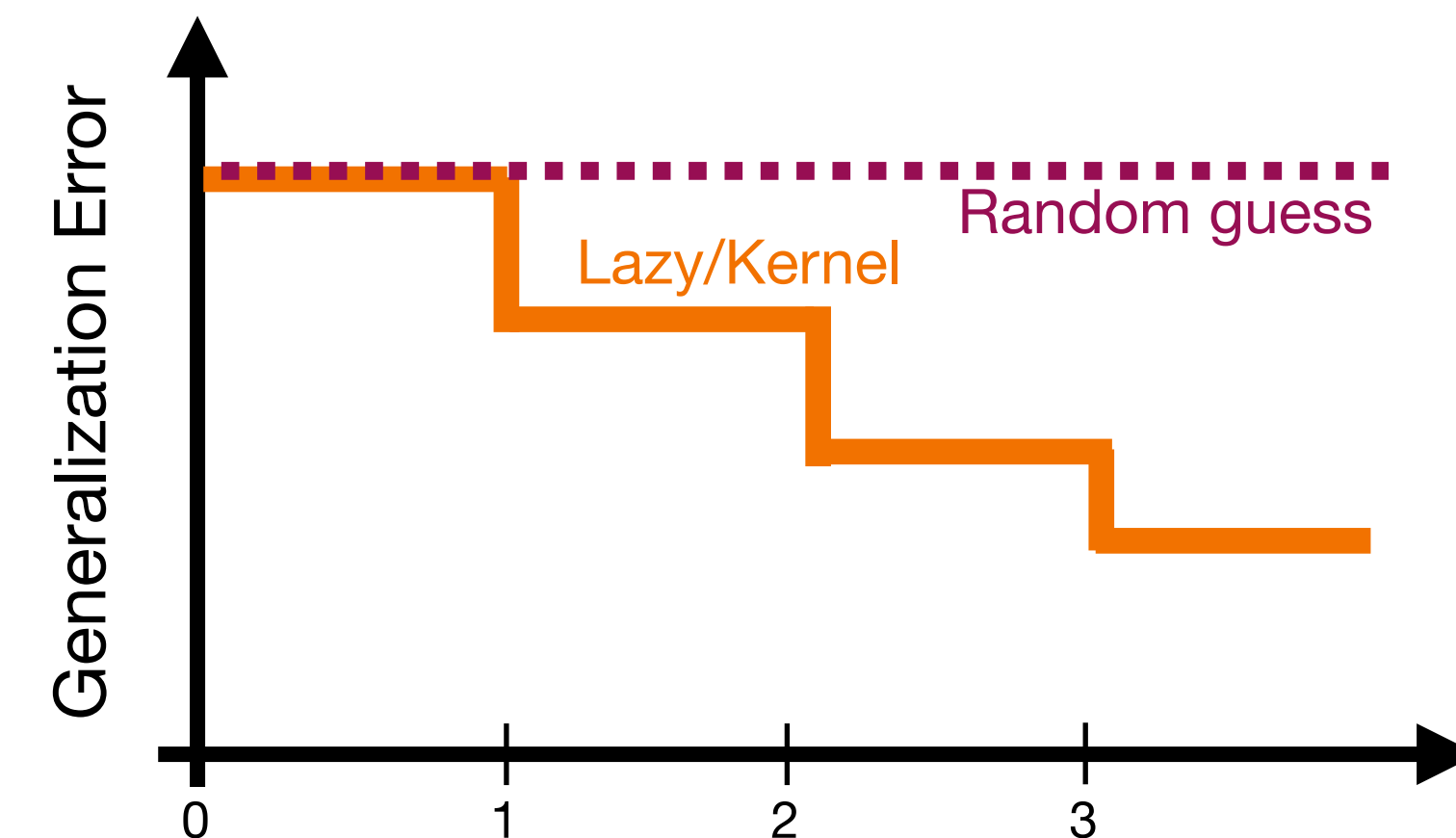
$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$

$$(n, p) = O(d)$$

$$(n, p) = O(d^2)$$

$$(n, p) = O(d^3)$$



See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]

# Recap of Sample Complexity

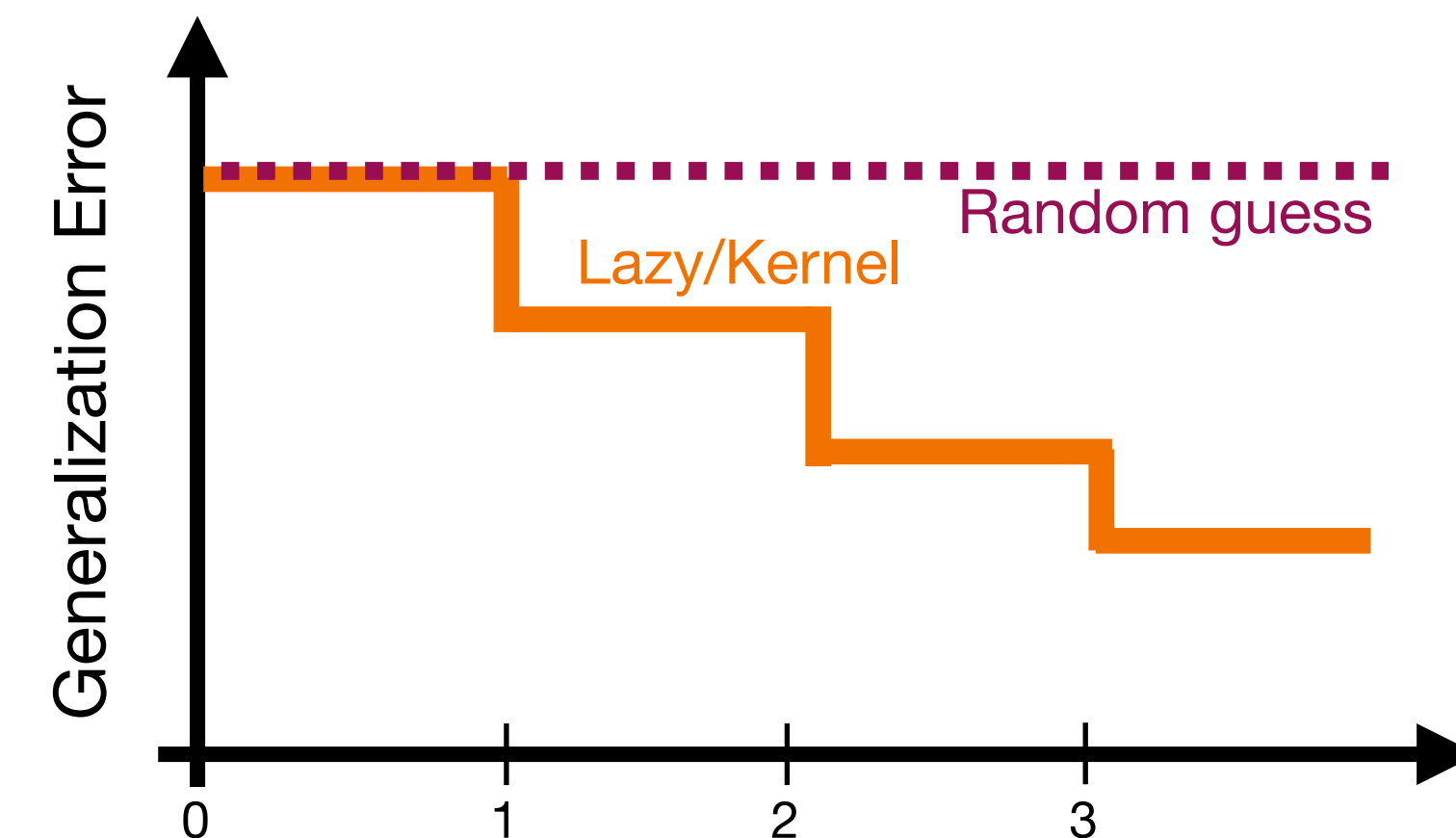
**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$
$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$

$(n, p) = O(d) \qquad (n, p) = O(d^2) \qquad (n, p) = O(d^3)$

**No adaptivity  $\implies$  searching in a  $\mathcal{O}(d^k)$  dimensional subspace of polynomials\***



See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]



# Recap of Sample Complexity

**Theorem (Informal)** [Mei, Misiakiewicz, Montanari '22]

In absence of feature learning (i.e. at initialization) one can only learn a **polynomial approximation of  $f^\star$  of degree  $\kappa$  with  $\min(n, p) = O(d^\kappa)$**

$$f^\star(\mathbf{x}) = \text{cst} + \sum_i \mu_i^{(1)} h_i^\star + \sum_{ij} \mu_{ij}^{(2)} h_i^\star h_j^\star + \sum_{ijk} \mu_{ijk}^{(3)} h_i^\star h_j^\star h_k^\star + \dots$$

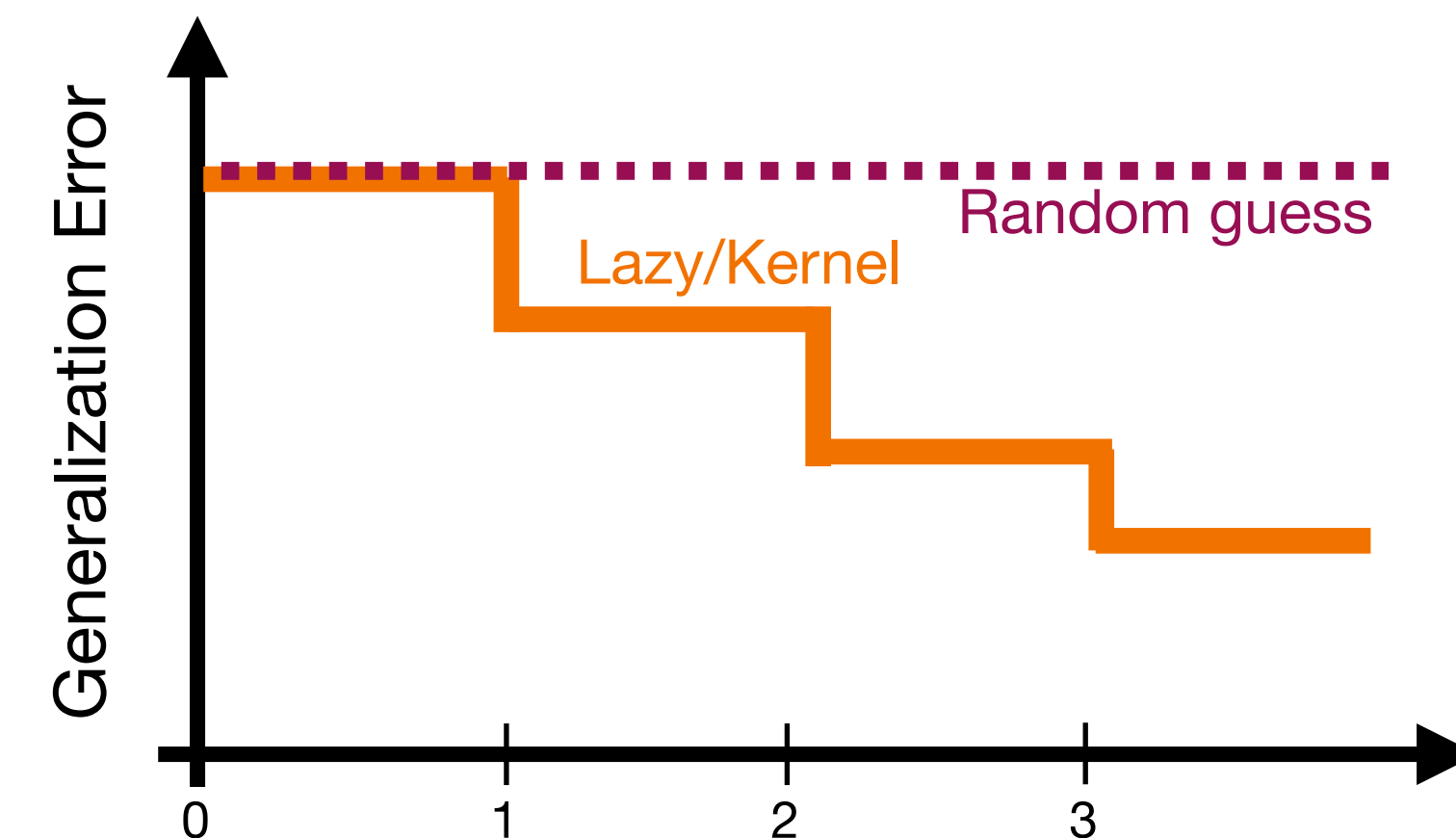
$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (f^\star(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x}))^2$$

$$(n, p) = O(d)$$

$$(n, p) = O(d^2)$$

$$(n, p) = O(d^3)$$

**No adaptivity  $\implies$  searching in a  $\mathcal{O}(d^k)$  dimensional subspace of polynomials\***



**\*possibly excluding a few “special” polynomials.**

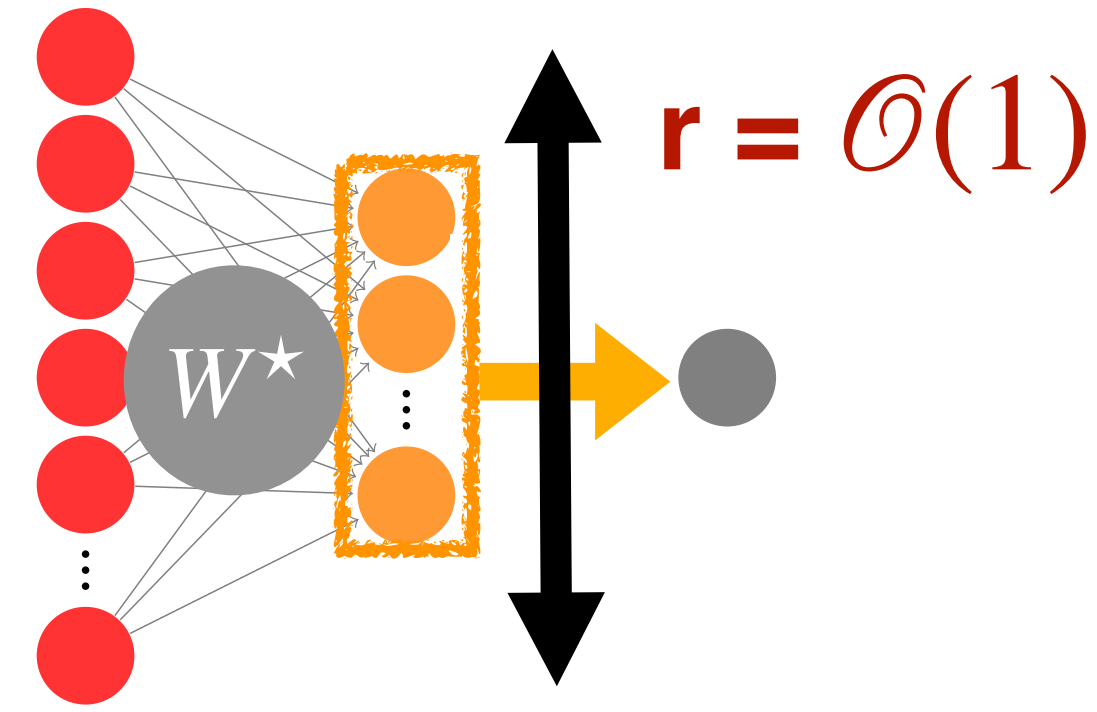
See also [El Karaoui '10; Mei-Montanari '19; Gerace '20; Jacot, Simsek, Spadaro, Hongler, Gabriel '20; Hu, Lu, '20; Dhifallah, Lu '20; Loureiro, Gerbelot, Cui, Goldt '21; Montanari & Saeed '22; Xiao, Hu, Misiakiewicz, Lu, Pennington '22; Dandi '23; Aguirre-López, Franz, Pastore '24]

# Advantage of Feature Learning



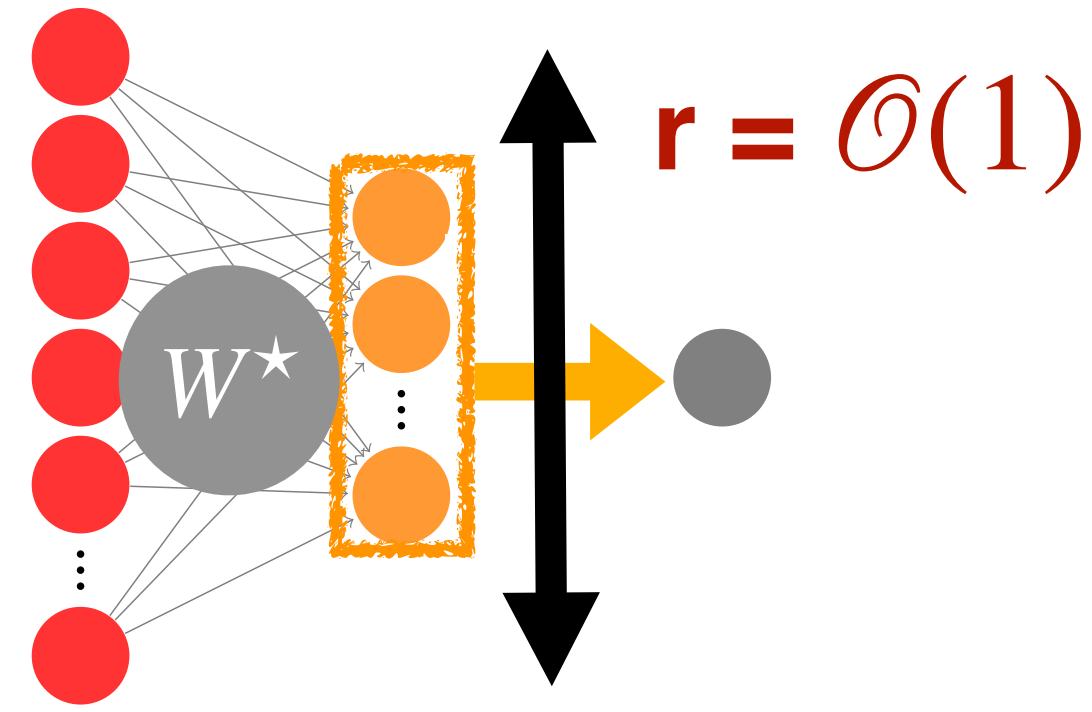
# Advantage of Feature Learning

$$y = f^{\star}(\mathbf{x}) = \mathbf{g}^{\star}(\mathbf{x}_{\star}) = \mathbf{W}^{\star} \mathbf{x}$$



# Advantage of Feature Learning

$$y = f^{\star}(\mathbf{x}) = \mathbf{g}^{\star}(\mathbf{x}_{\star}) = \mathbf{W}^{\star} \mathbf{x}$$

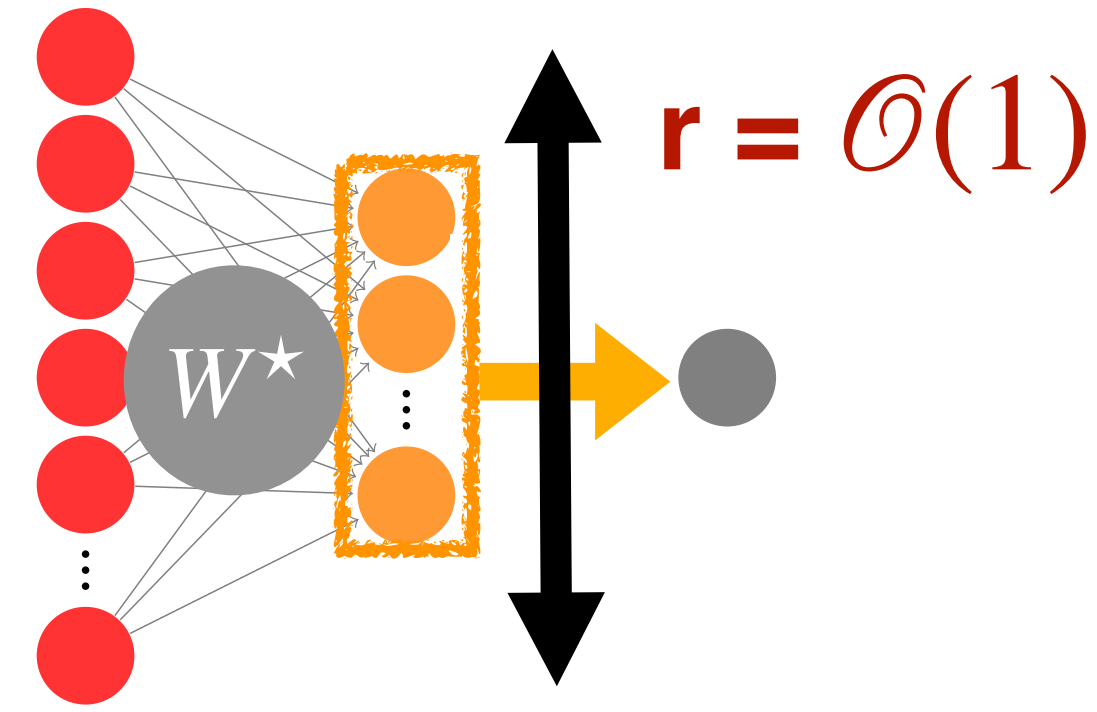


**Information/  
generative exponent  
 $\leq 2$**



# Advantage of Feature Learning

$$y = f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x})$$

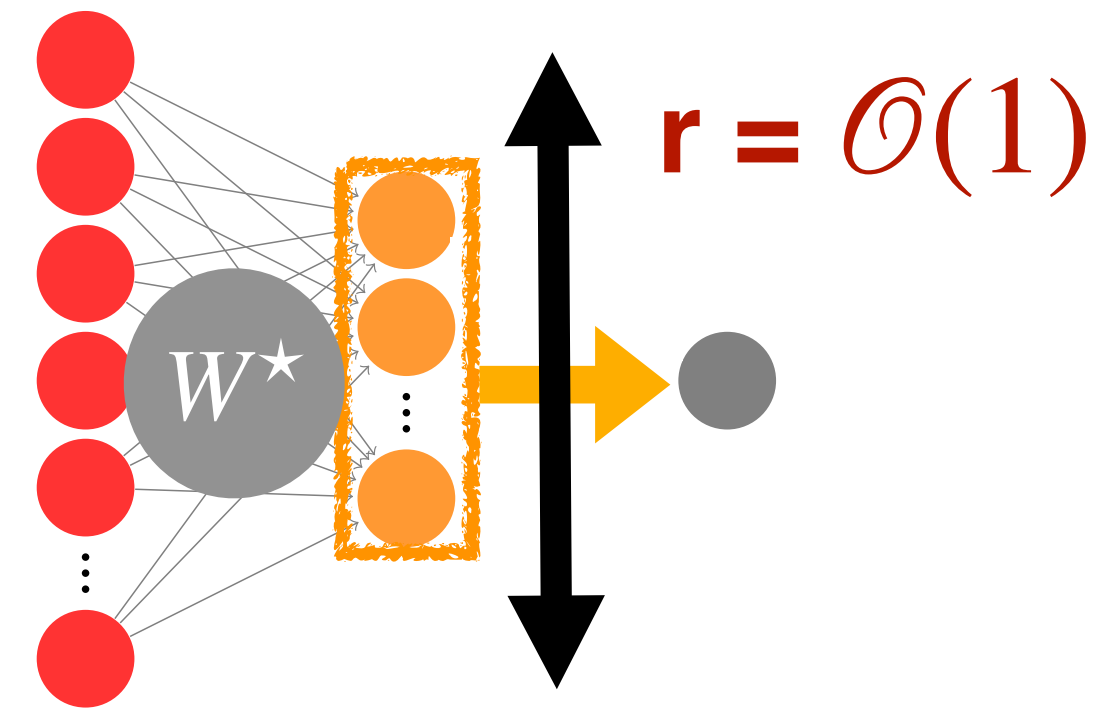


**Information/  
generative exponent  
 $\leq 2$**

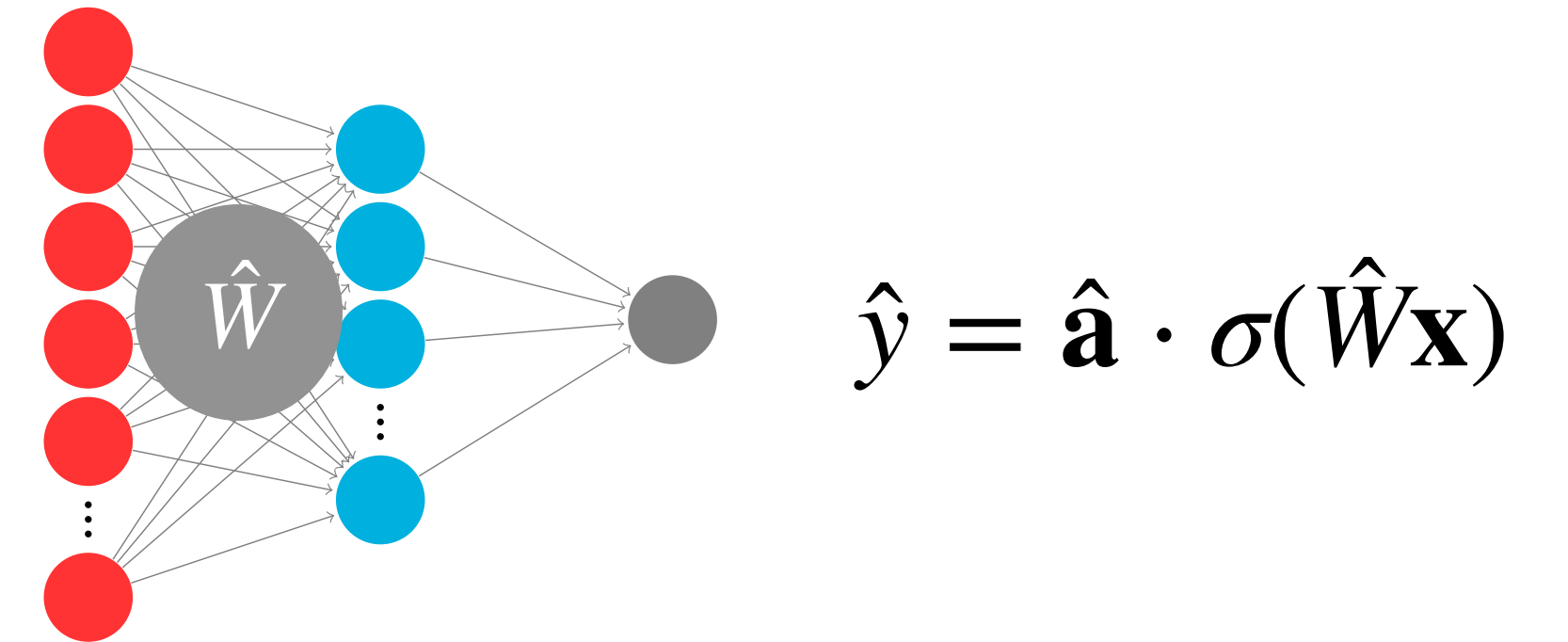


# Advantage of Feature Learning

$$y = f^{\star}(\mathbf{x}) = \mathbf{g}^{\star}(\mathbf{x}_{\star}) = \mathbf{W}^{\star} \mathbf{x}$$

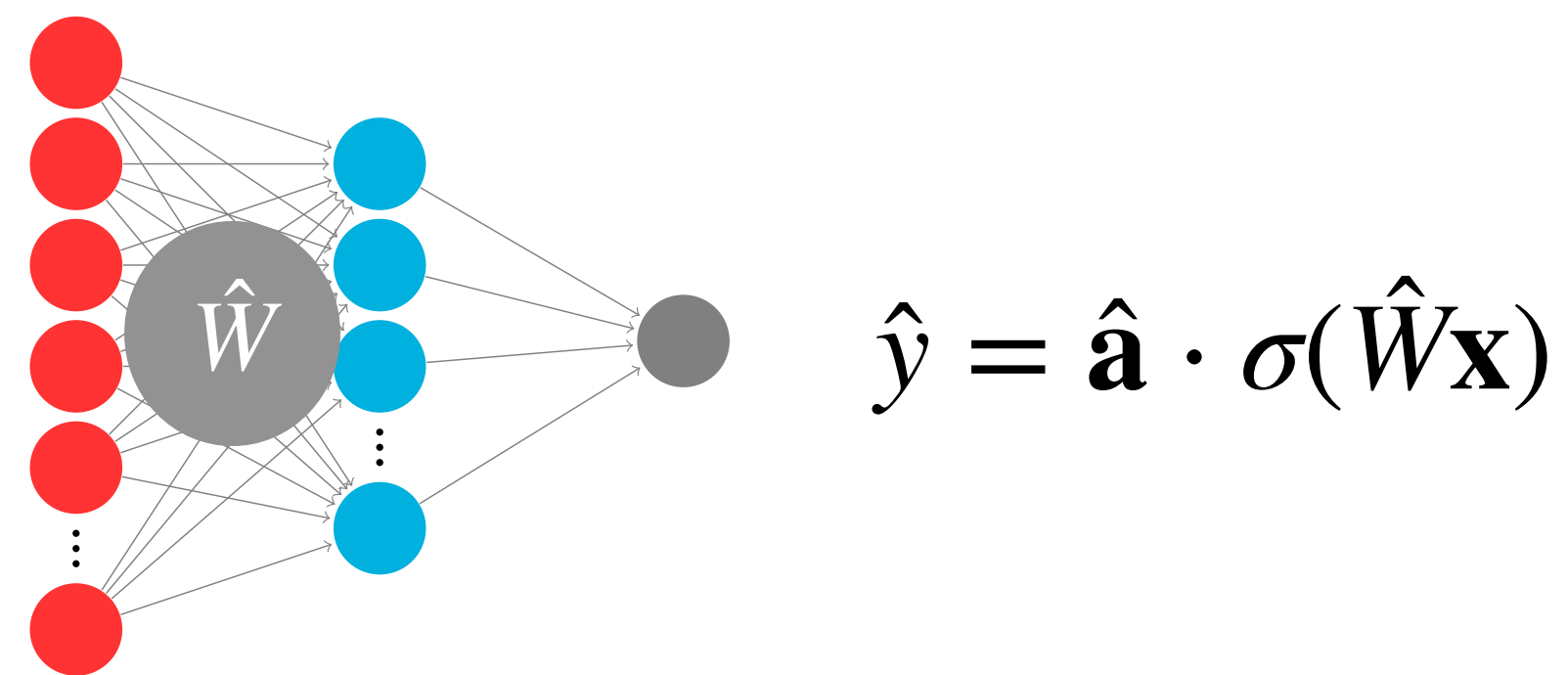
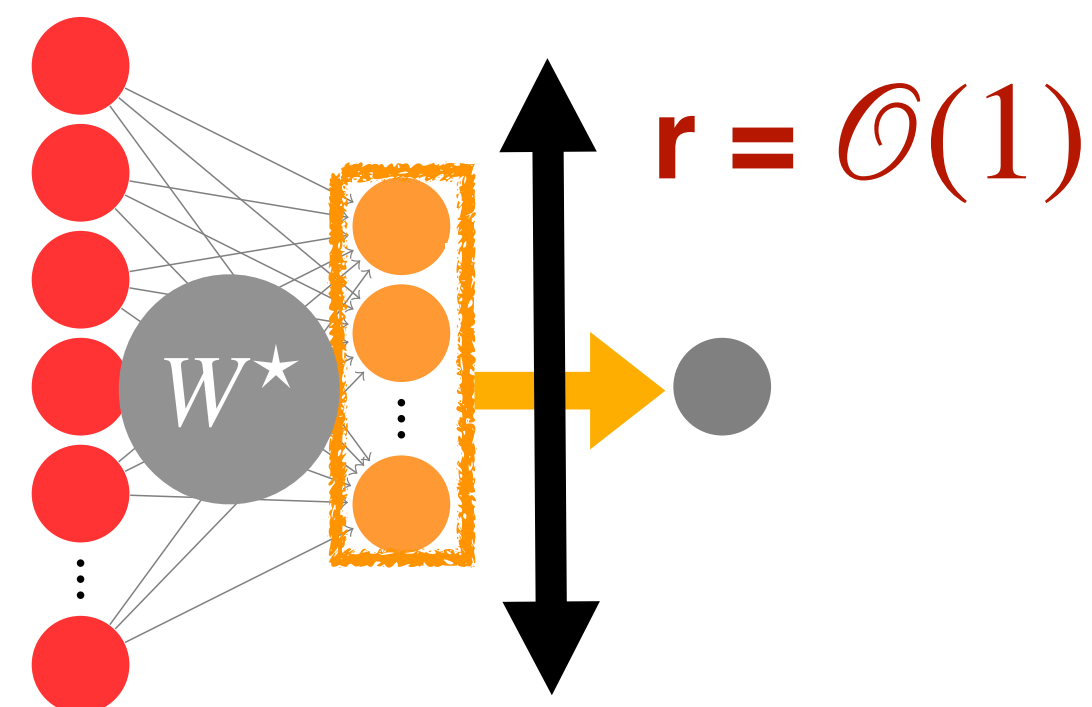


**Information/  
generative exponent  
 $\leq 2$**



# Advantage of Feature Learning

$$y = f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x})$$



With  $\tilde{O}(d)$  samples, GD/SGD on  $\hat{W}$  recovers

**Information/  
generative exponent  
 $\leq 2$**

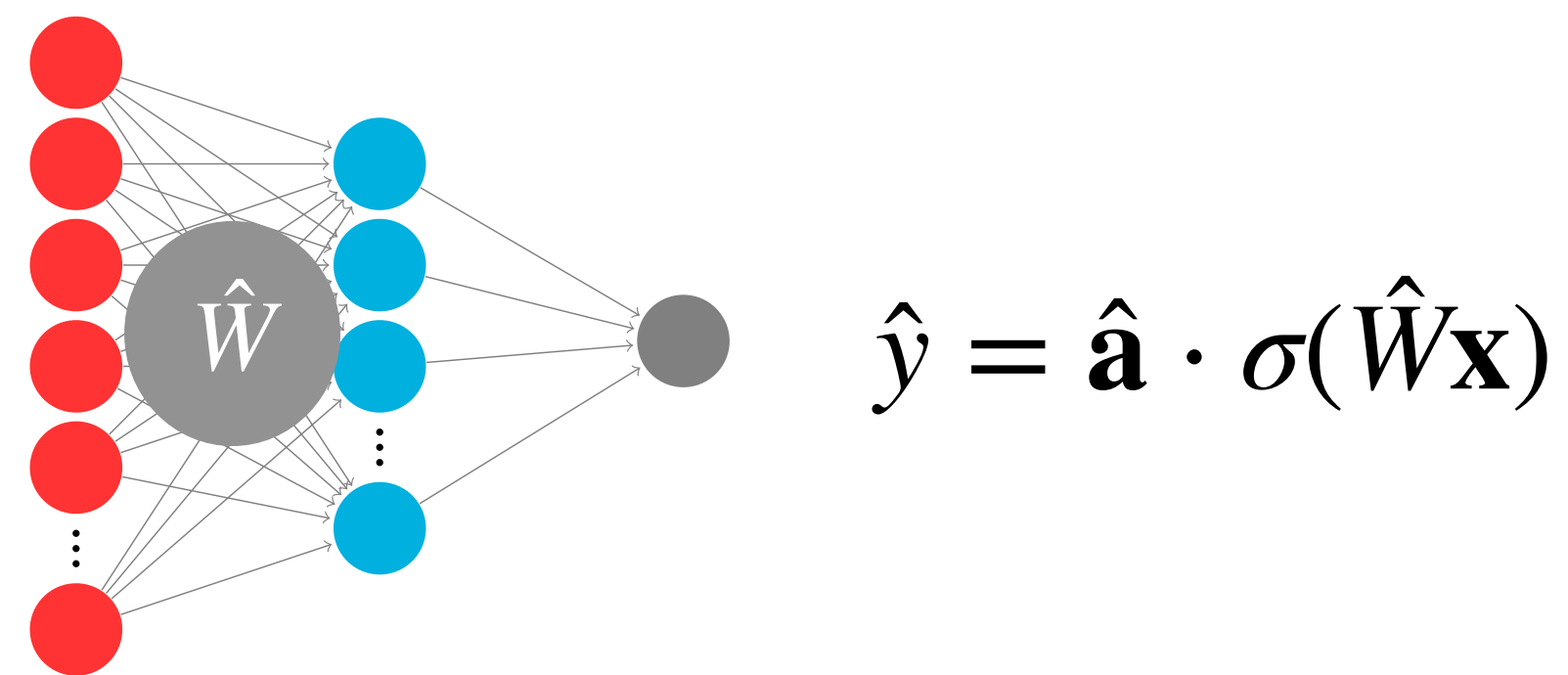
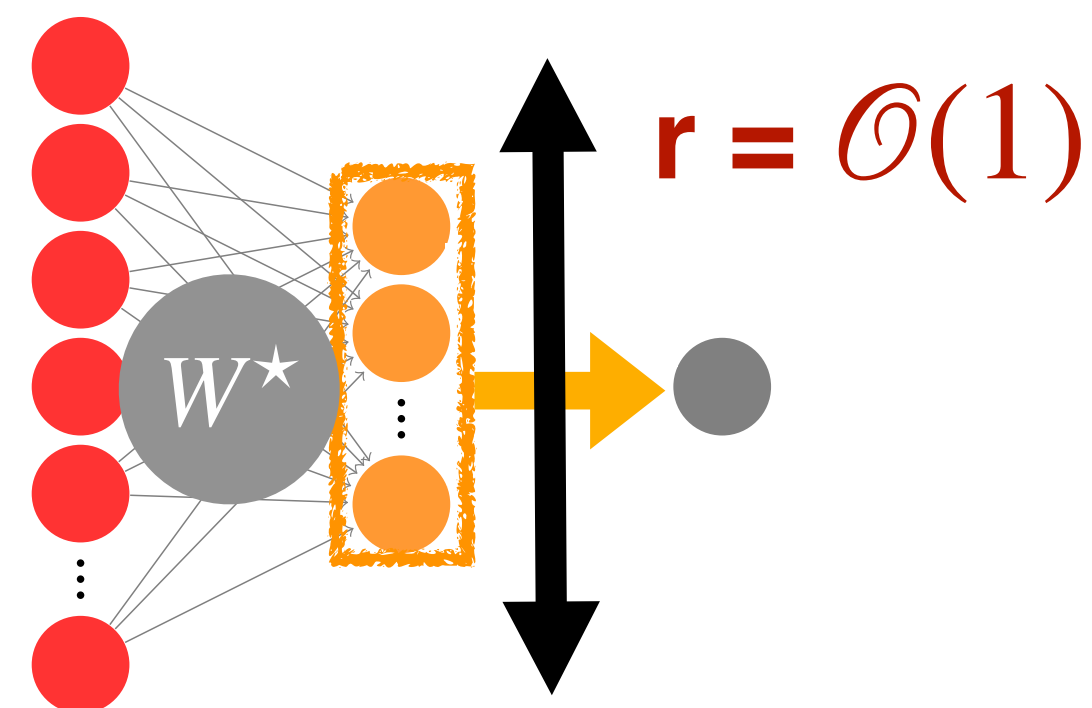


$$\hat{W} \approx Z_1 W^\star + Z_2$$



# Advantage of Feature Learning

$$y = f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x})$$



**Information/  
generative exponent  
 $\leq 2$**

With  $\tilde{O}(d)$  samples, GD/SGD on  $\hat{W}$  recovers

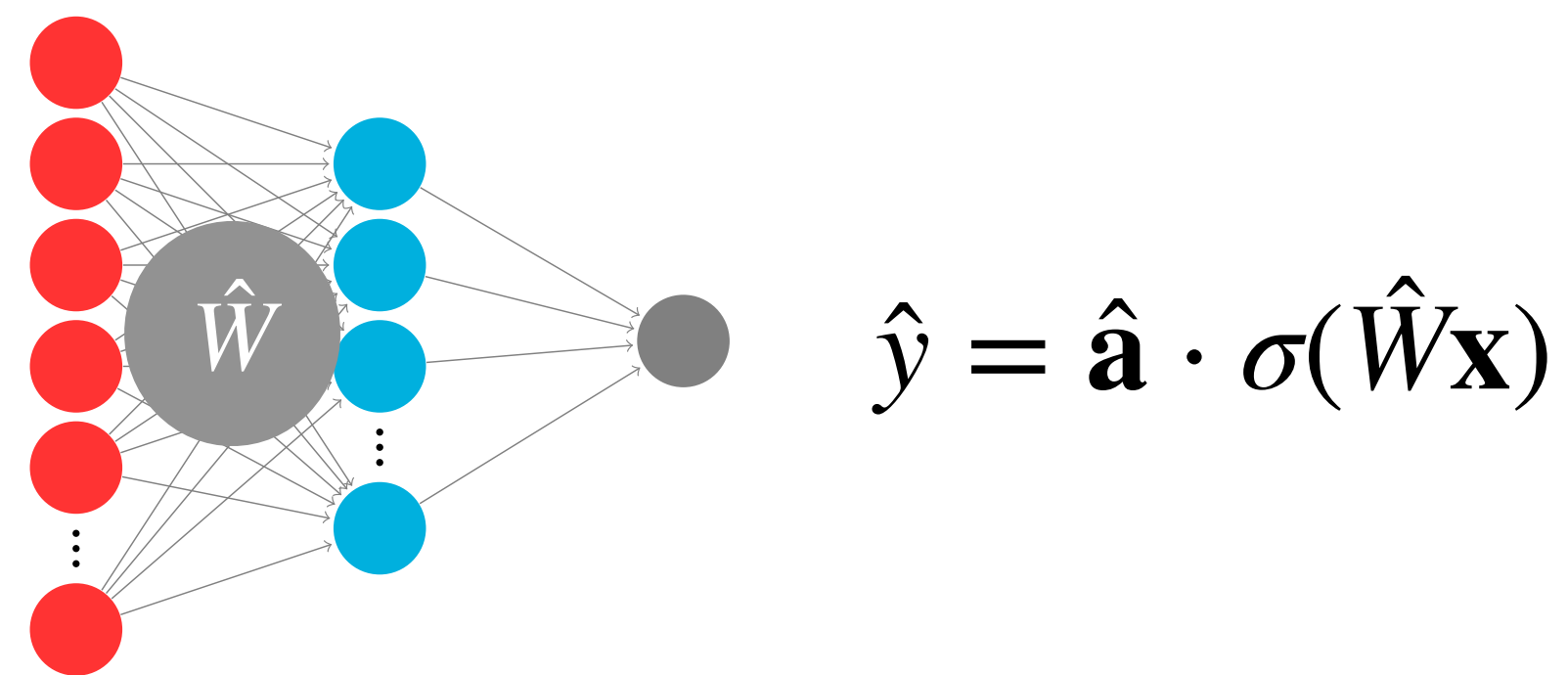
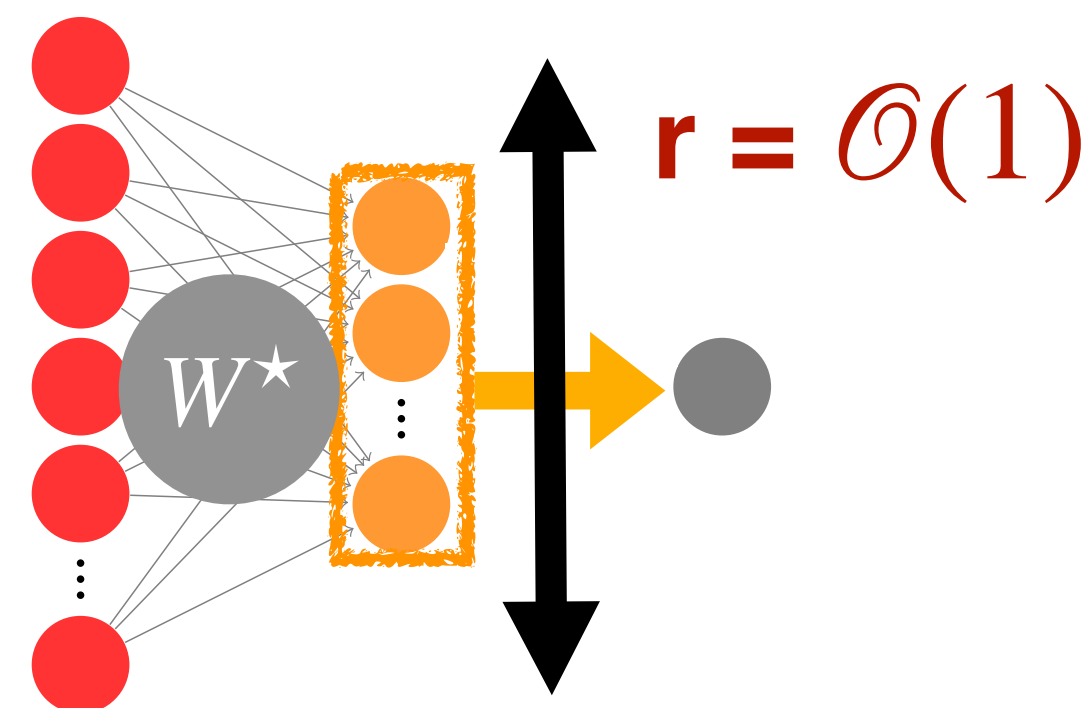
$$\hat{W} \approx Z_1 W^\star + Z_2$$

Then the neural net is equivalent to

$$\hat{y} \approx \hat{\mathbf{a}} \cdot \sigma((Z_1 W^\star + Z_2)\mathbf{x}) = \hat{\mathbf{a}} \cdot \sigma(\mathbf{Z}_1 \mathbf{x}^\star + \mathbf{Z}_3)$$

# Advantage of Feature Learning

$$y = f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x})$$



**Information/  
generative exponent  
 $\leq 2$**

With  $\tilde{O}(d)$  samples, GD/SGD on  $\hat{W}$  recovers

$$\hat{W} \approx Z_1 W^\star + Z_2$$

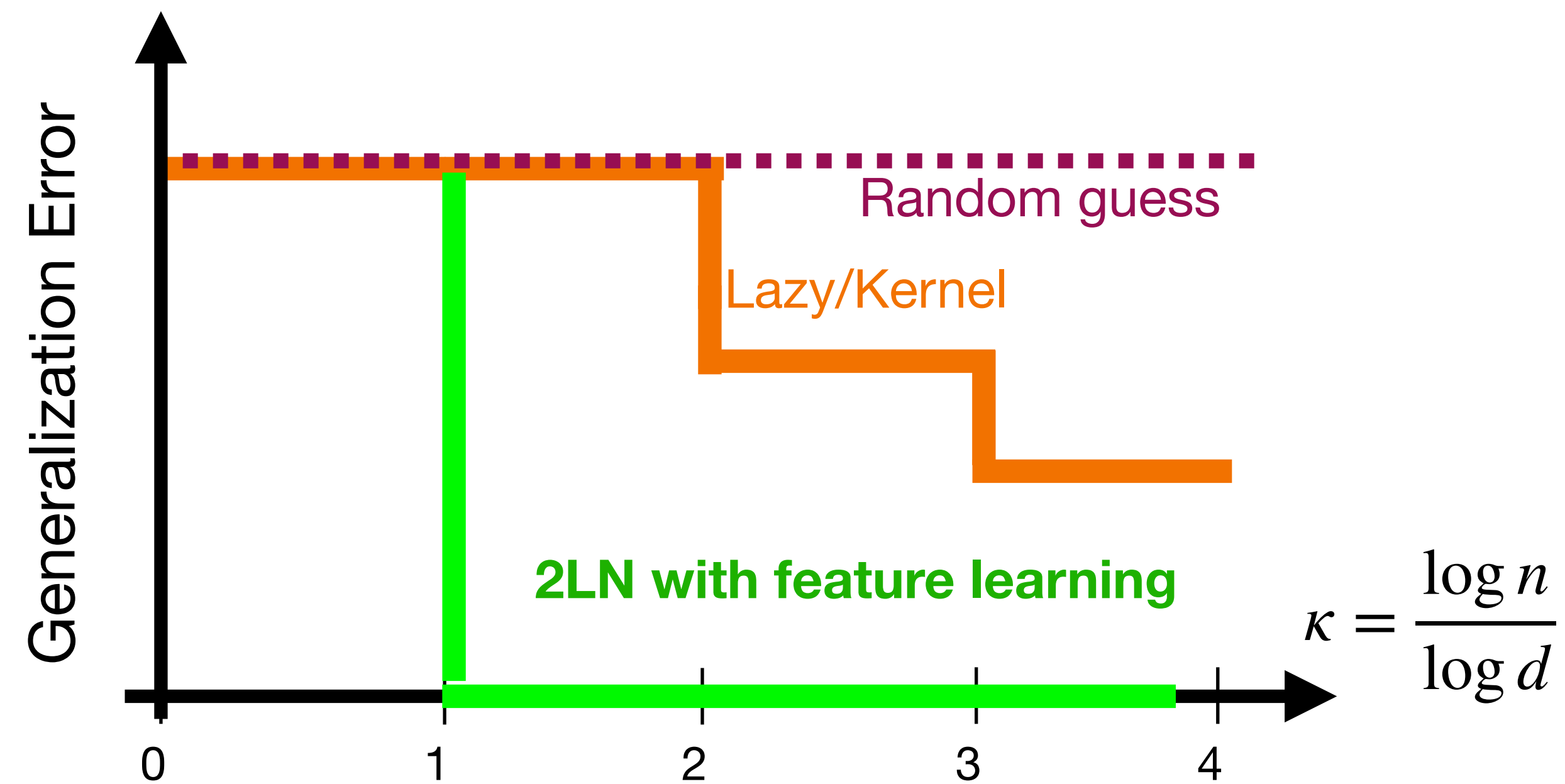
Then the neural net is equivalent to

$$\hat{y} \approx \hat{\mathbf{a}} \cdot \sigma((Z_1 W^\star + Z_2)\mathbf{x}) = \hat{\mathbf{a}} \cdot \sigma(\mathbf{Z}_1 \mathbf{x}^\star + \mathbf{Z}_3)$$

Random feature in (finite)  
reduced space

$$d \rightarrow d^{\text{eff}} = r = O(1)$$

# Dimensionality reduction



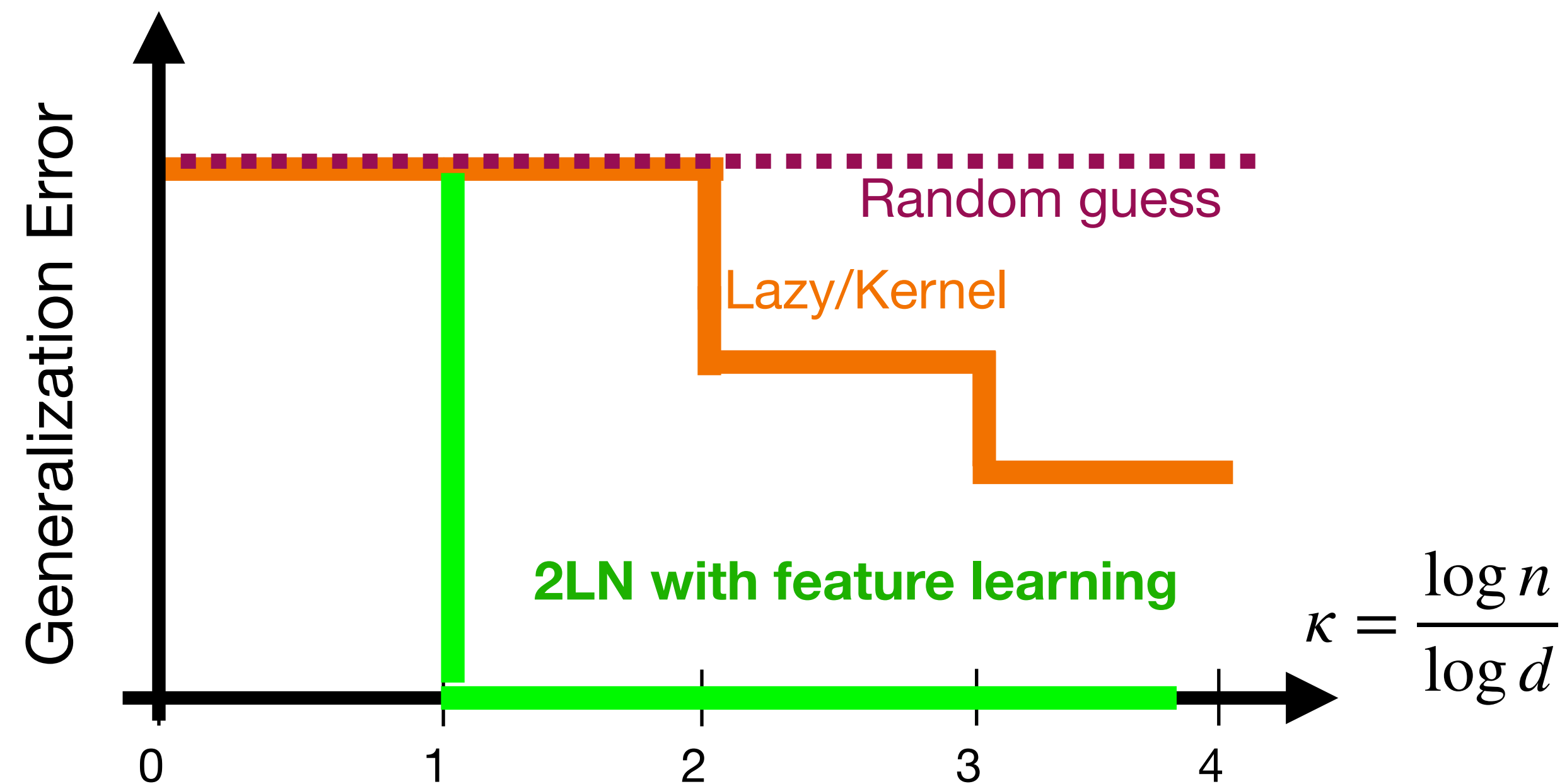
**Sample Complexity Reduction**

$$\mathcal{O}(d^k) \rightarrow \tilde{\mathcal{O}}(d)$$



# Dimensionality reduction

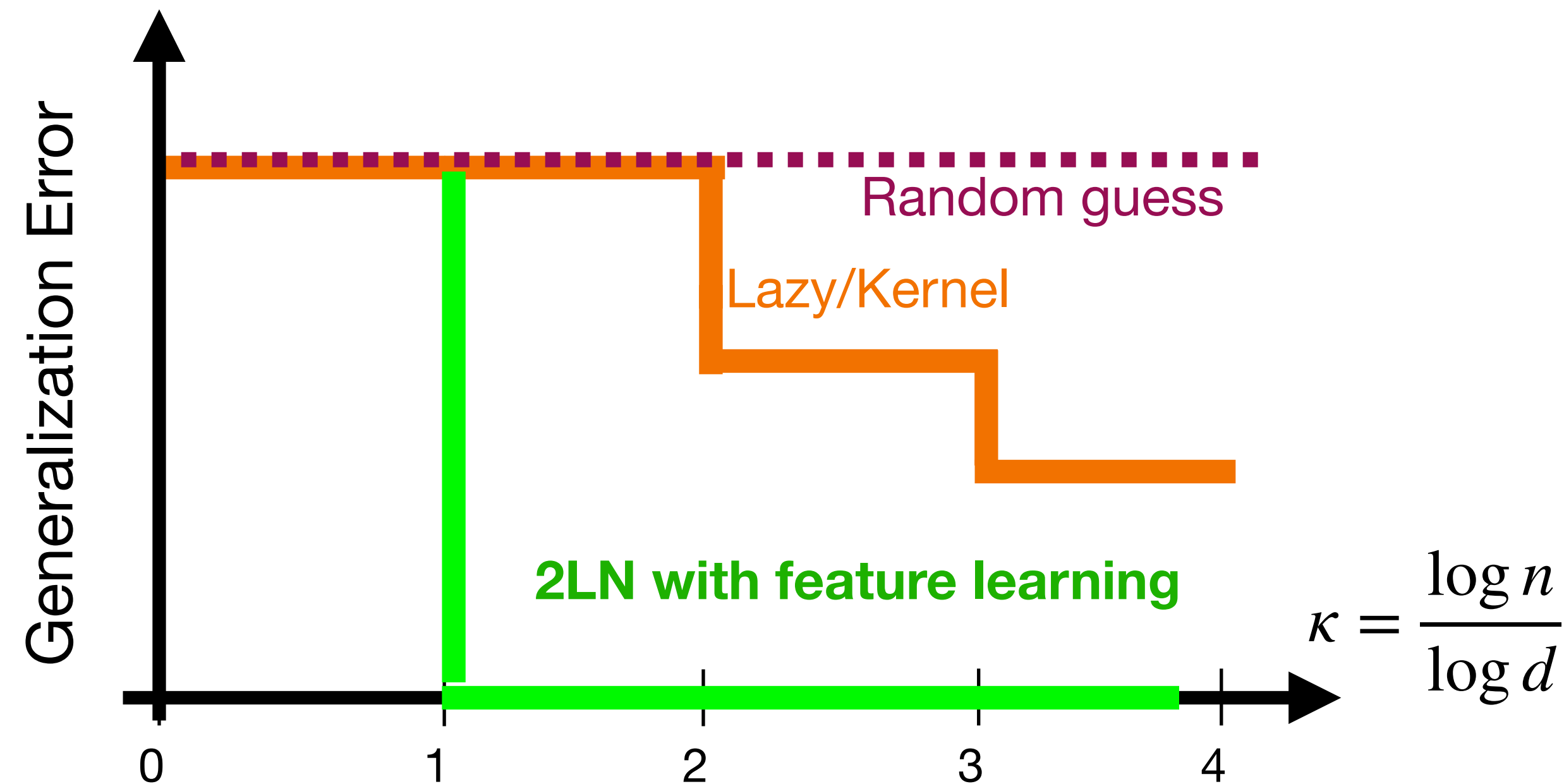
searching in a  $\mathcal{O}(d^k)$   
dimensional subspace of  
polynomials



Sample Complexity Reduction  
 $\mathcal{O}(d^k) \rightarrow \tilde{\mathcal{O}}(d)$

# Dimensionality reduction

searching in a  $\mathcal{O}(d^k)$   
dimensional subspace of  
polynomials



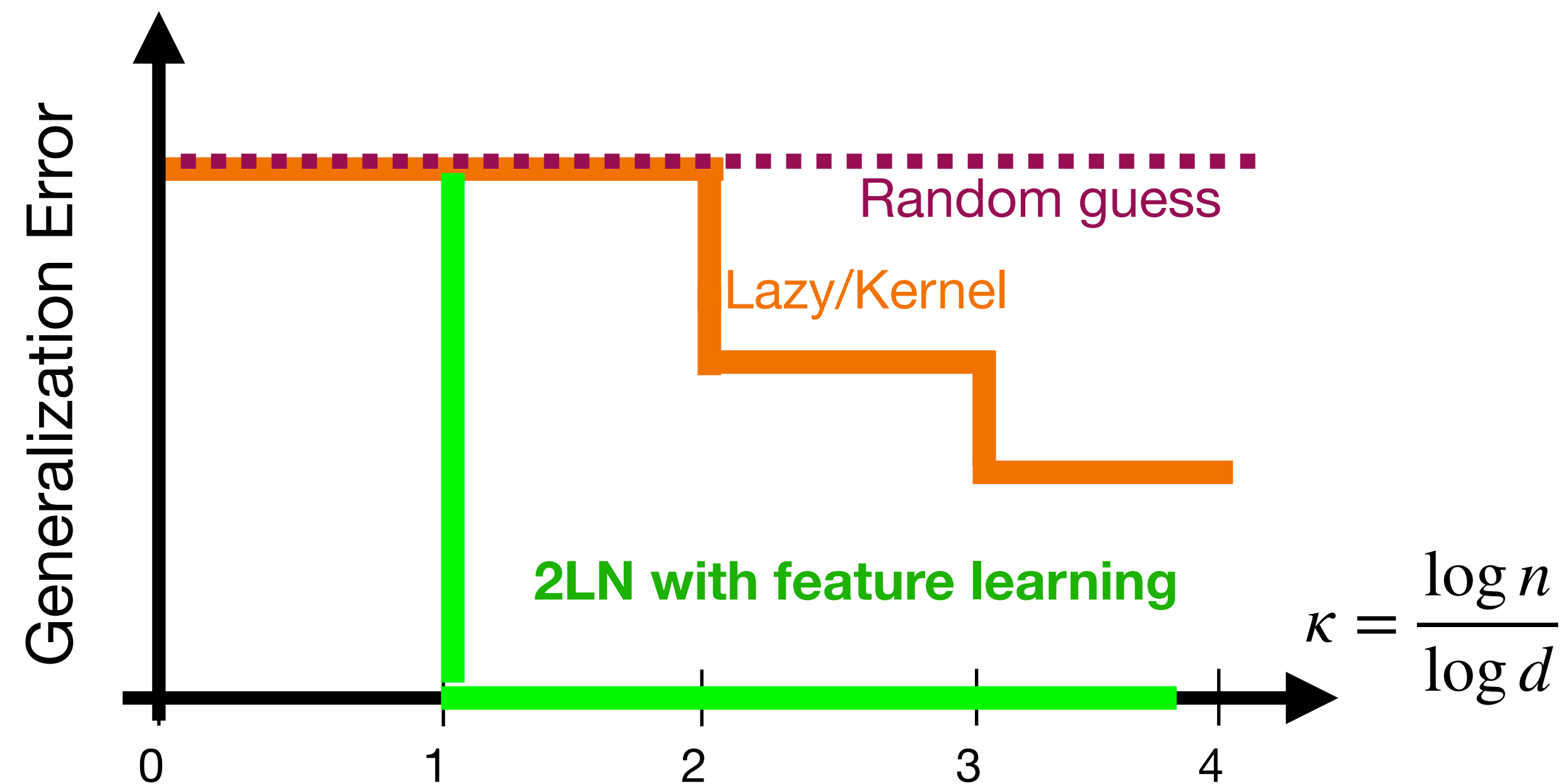
Sample Complexity Reduction  
 $\mathcal{O}(d^k) \rightarrow \tilde{\mathcal{O}}(d)$

# Dimensionality reduction

searching in a  $\mathcal{O}(d^k)$   
dimensional subspace of  
polynomials



searching in a  $\mathcal{O}(r^k)$   
dimensional subspace of  
polynomials



Sample Complexity Reduction  
 $\mathcal{O}(d^k) \rightarrow \tilde{\mathcal{O}}(d)$

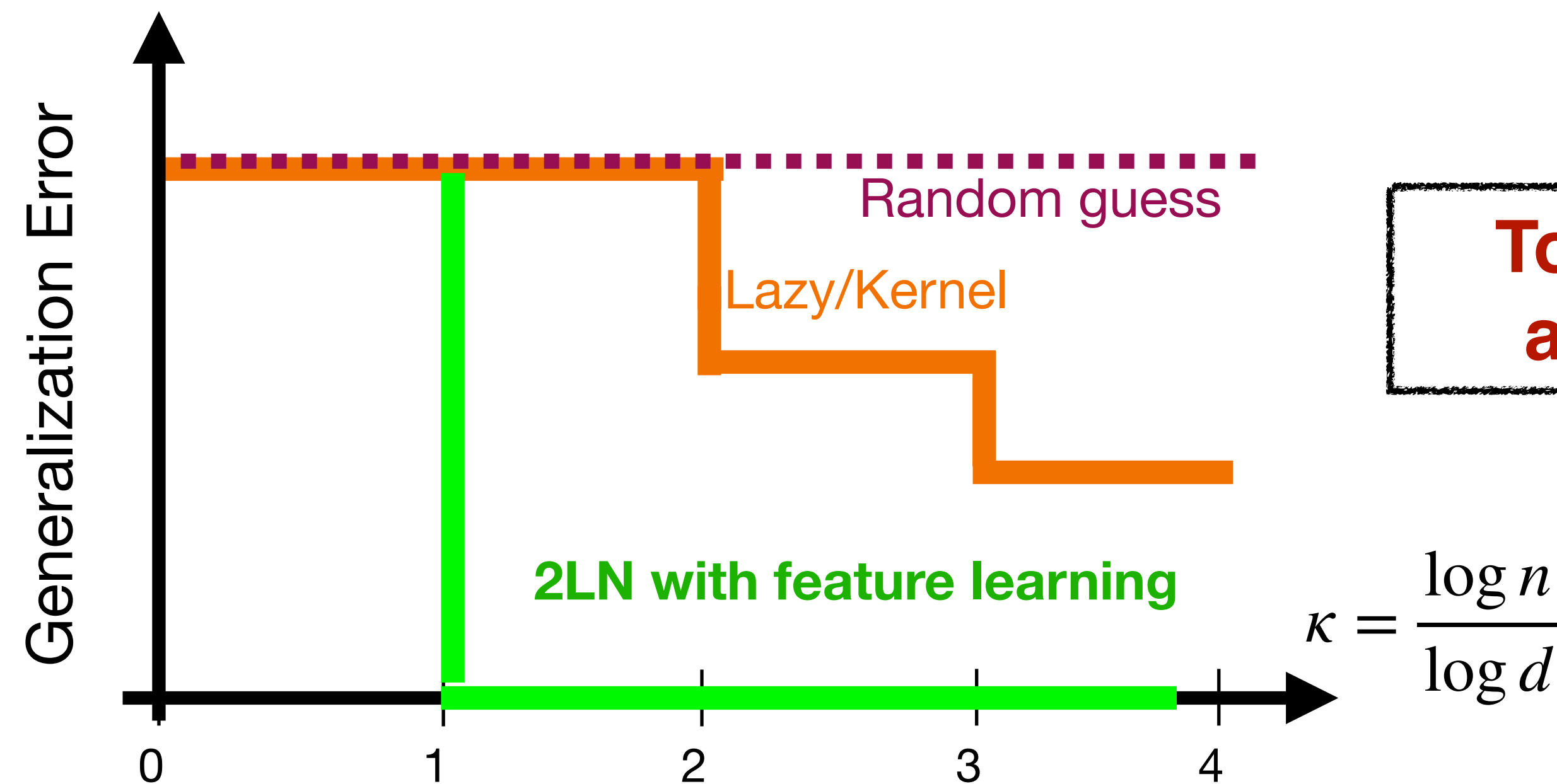


# Dimensionality reduction

searching in a  $\mathcal{O}(d^k)$   
dimensional subspace of  
polynomials



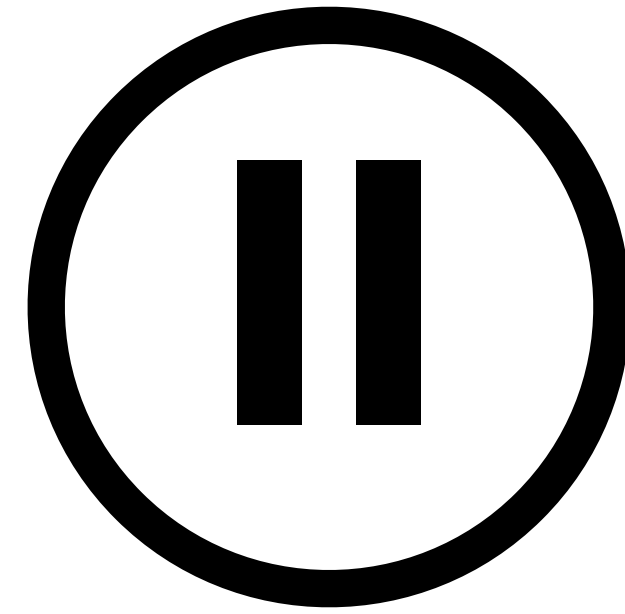
searching in a  $\mathcal{O}(r^k)$   
dimensional subspace of  
polynomials



Too easy...no further  
advantage of depth

Sample Complexity Reduction  
 $\mathcal{O}(d^k) \rightarrow \tilde{\mathcal{O}}(d)$





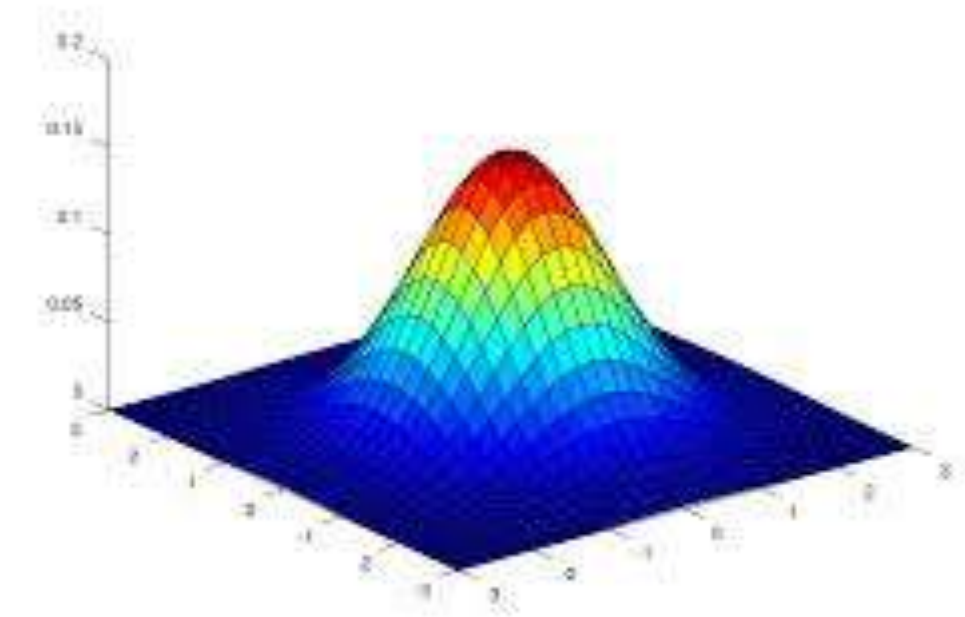
**Can we generalize this picture to  
arbitrary depth?**



# A SIGHT beyond Two Layers (Single-Index Gaussian Hierarchical Targets)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\varepsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$P_k$ : polynomial of degree  $k$   
 $g^\star$ : non-linearity

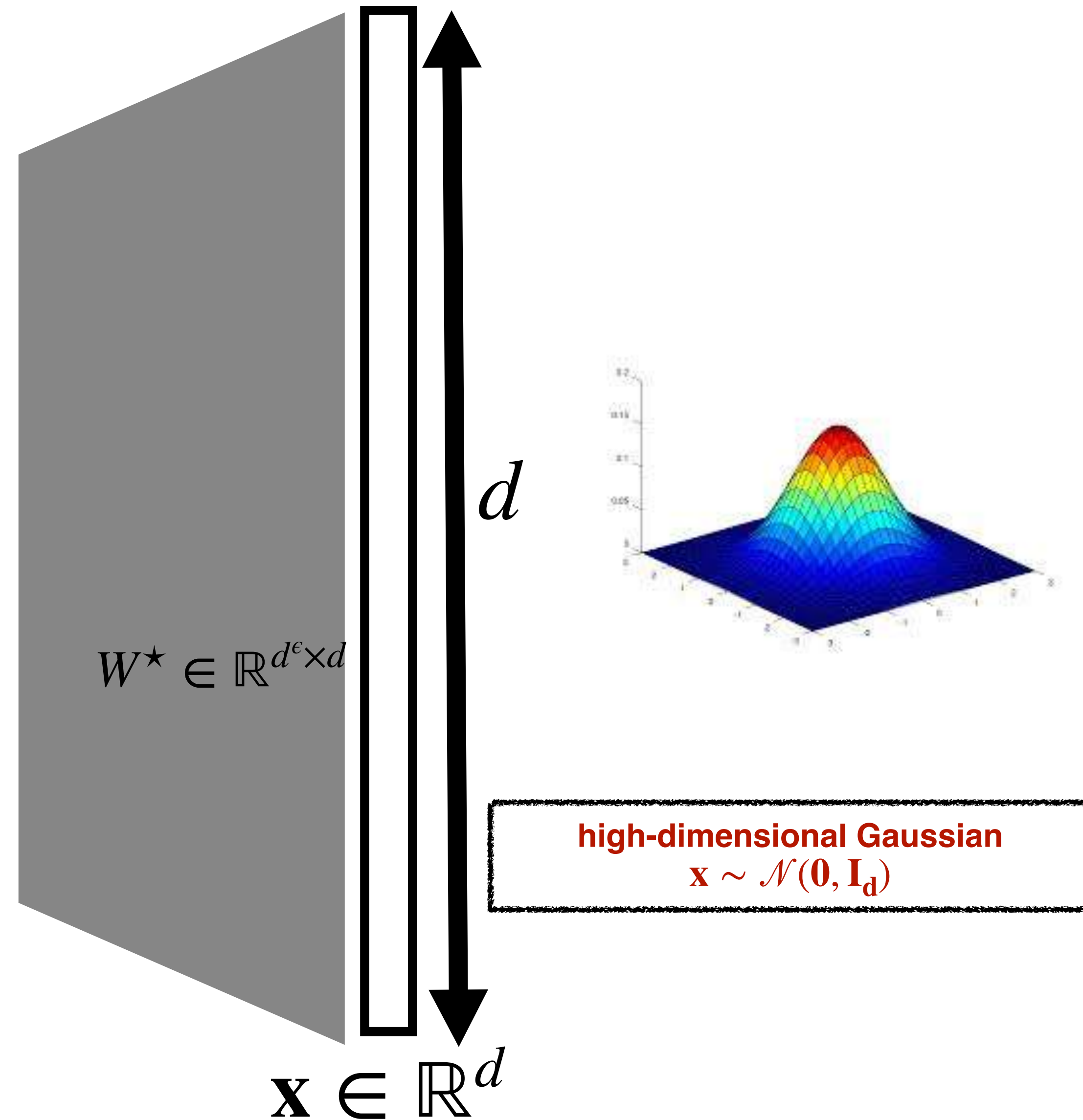


high-dimensional Gaussian  
 $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$

# A SIGHT beyond Two Layers (Single-Index Gaussian Hierarchical Targets)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$P_k$ : polynomial of degree  $k$   
 $g^\star$ : non-linearity

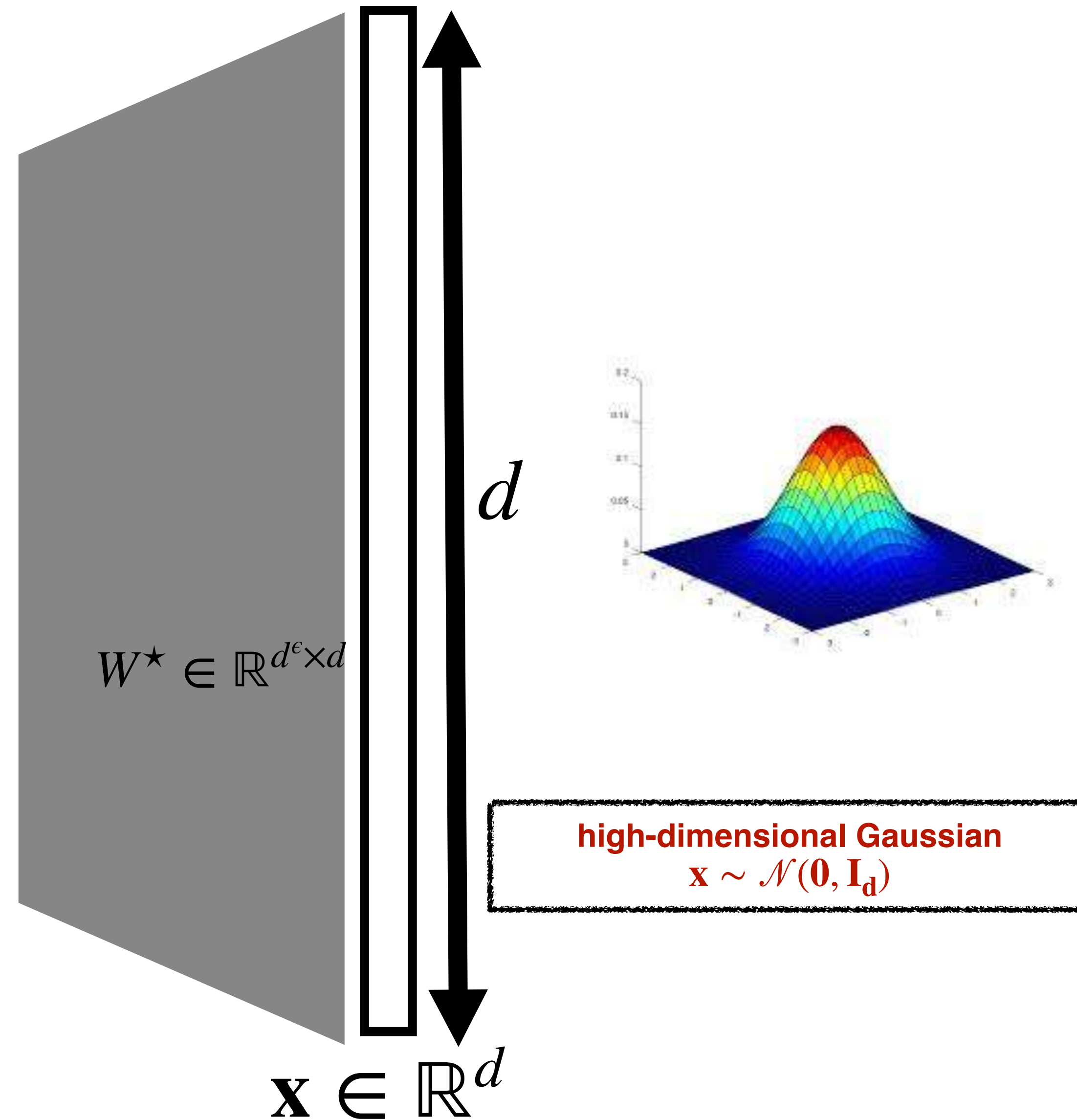


# A SIGHT beyond Two Layers (Single-Index Gaussian Hierarchical Targets)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$P_k$ : polynomial of degree  $k$   
 $g^\star$ : non-linearity

Subspace dimension  $d^\epsilon$   
for  $0 < \epsilon < 1$



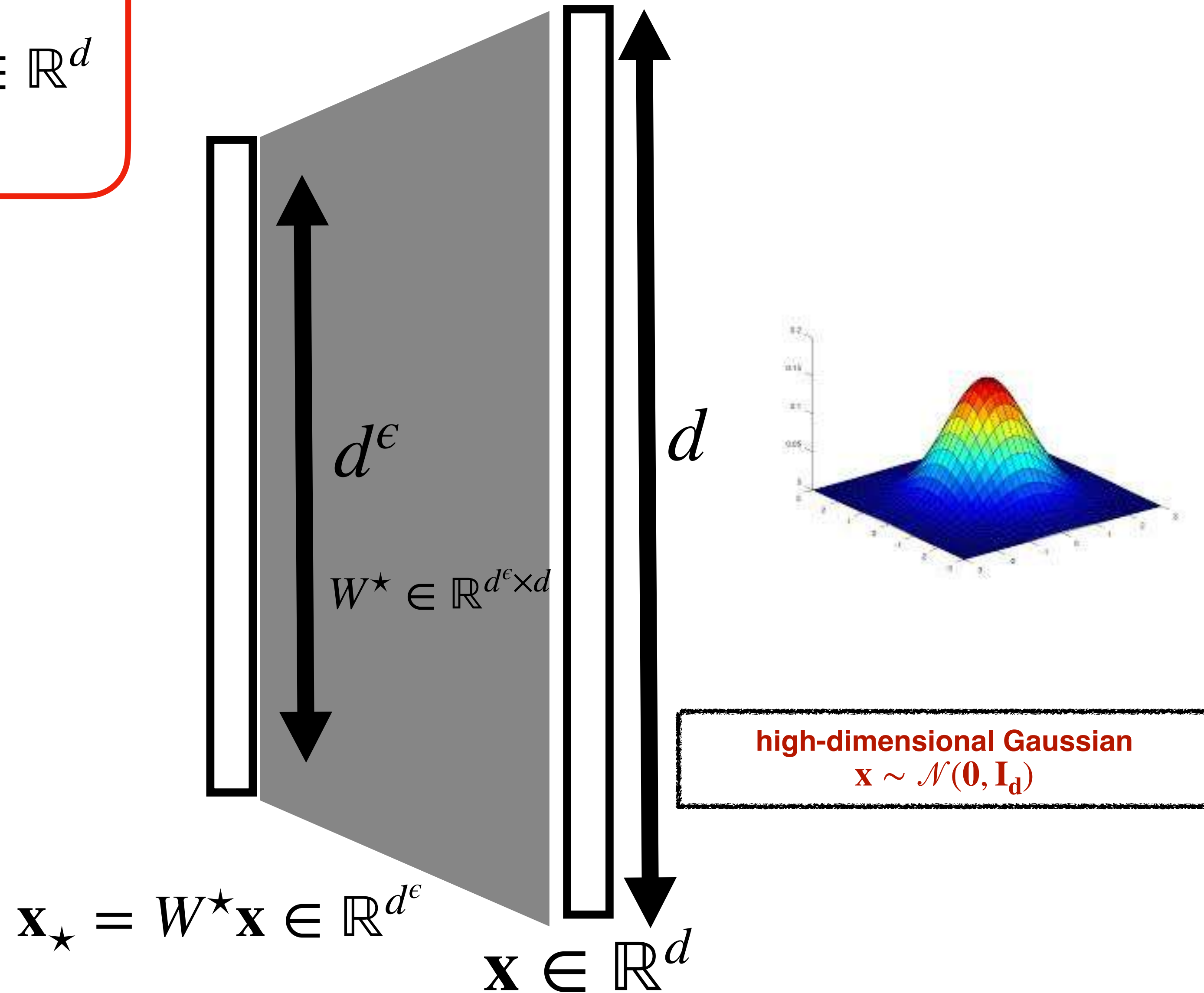


# A SIGHT beyond Two Layers (Single-Index Gaussian Hierarchical Targets)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$P_k$ : polynomial of degree  $k$   
 $g^\star$ : non-linearity

Subspace dimension  $d^\epsilon$   
for  $0 < \epsilon < 1$



# A SIGHT beyond Two Layers (Single-Index Gaussian Hierarchical Targets)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$P_k$ : polynomial of degree  $k$   
 $g^\star$ : non-linearity

$$f^\star(\mathbf{x}) = g^\star(h^\star) \quad \square \text{---} \square \quad \mathbf{a}^\star$$

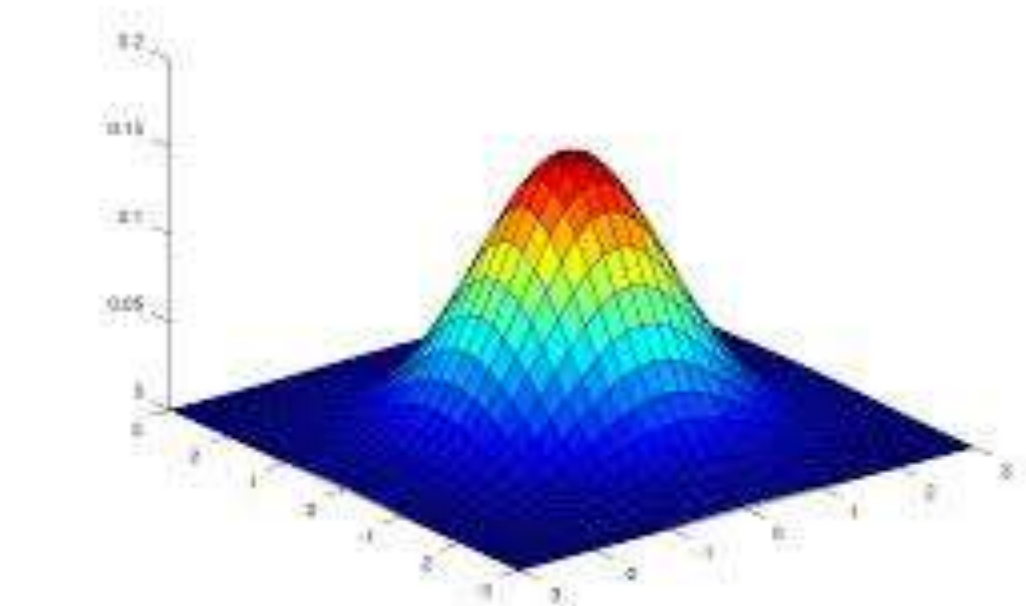
$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

Subspace dimension  $d^\epsilon$   
 for  $0 < \epsilon < 1$

$$\mathbf{x}_\star = W^\star \mathbf{x} \in \mathbb{R}^{d^\epsilon}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$W^\star \in \mathbb{R}^{d^\epsilon \times d}$$

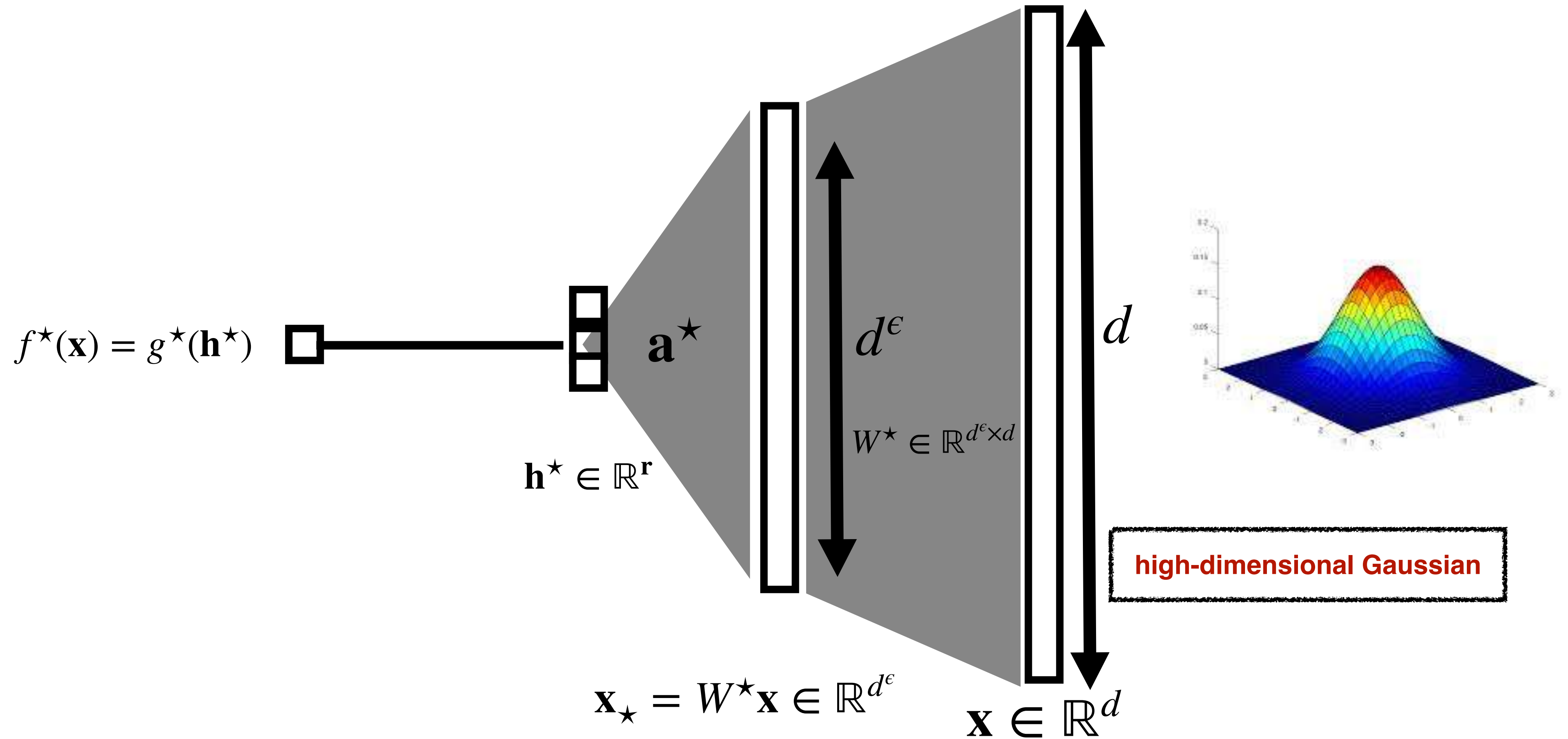


high-dimensional Gaussian  
 $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$

# MIGHT (Multi-Index Gaussian Hierarchical Targets)

$$f^{\star}(\mathbf{x}) = g^{\star}(\mathbf{h}_1^{\star}(\mathbf{x}), \dots, \mathbf{h}_r^{\star}(\mathbf{x})), \mathbf{x} \in \mathbb{R}^d$$

$$h_m^{\star}(\mathbf{x}) = \frac{1}{\sqrt{d^{\epsilon}}} \mathbf{a}_m^{\star \top} \mathbf{P}_{k,m} (\mathbf{W}_m^{\star} \mathbf{x}), m = 1, \dots, r$$





# Key features of SIGHTS/MIGHTS

# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.

# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .



# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .
- By CLT,  $h^*$  asymptotically Gaussian and independent.

# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .
- By CLT,  $h^*$  asymptotically Gaussian and independent.
- Motivated by Nichani et al. 2023:

# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .
- By CLT,  $h^*$  asymptotically Gaussian and independent.
- Motivated by Nichani et al. 2023:

$$f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}^\top \mathbf{A} \mathbf{x})$$



# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .
- By CLT,  $h^*$  asymptotically Gaussian and independent.
- Motivated by Nichani et al. 2023:

$$f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}^\top \mathbf{A} \mathbf{x}) \longrightarrow$$

# Key features of SIGHTS/MIGHTS

- SIGHT/MIGHT are simply three-layer neural networks.
- Dimension of features reduces from  $d$  to  $d^\epsilon$  to  $\mathcal{O}(1)$ .
- By CLT,  $h^*$  asymptotically Gaussian and independent.
- Motivated by Nichani et al. 2023:

$$f^\star(\mathbf{x}) = \mathbf{g}^\star(\mathbf{x}^\top \mathbf{A} \mathbf{x}) \longrightarrow$$

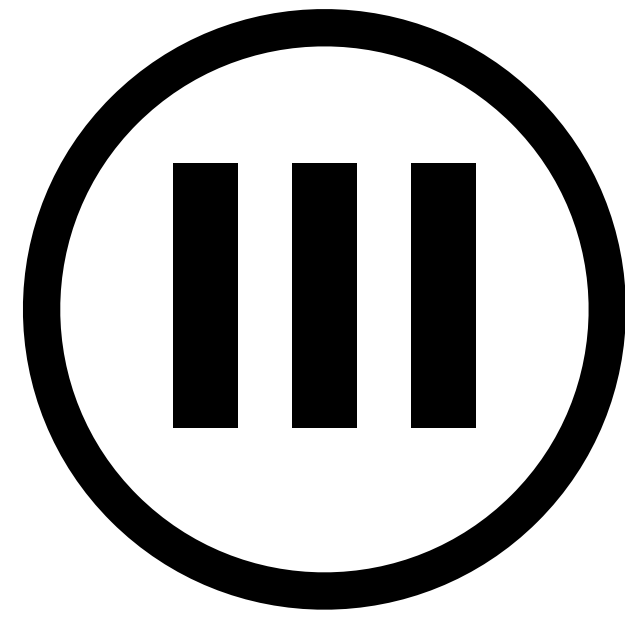
**2-layer NN can't learn  
non-linear features**

**1 level of dimension-  
reduction**

**No adaptivity of first layer  
required**



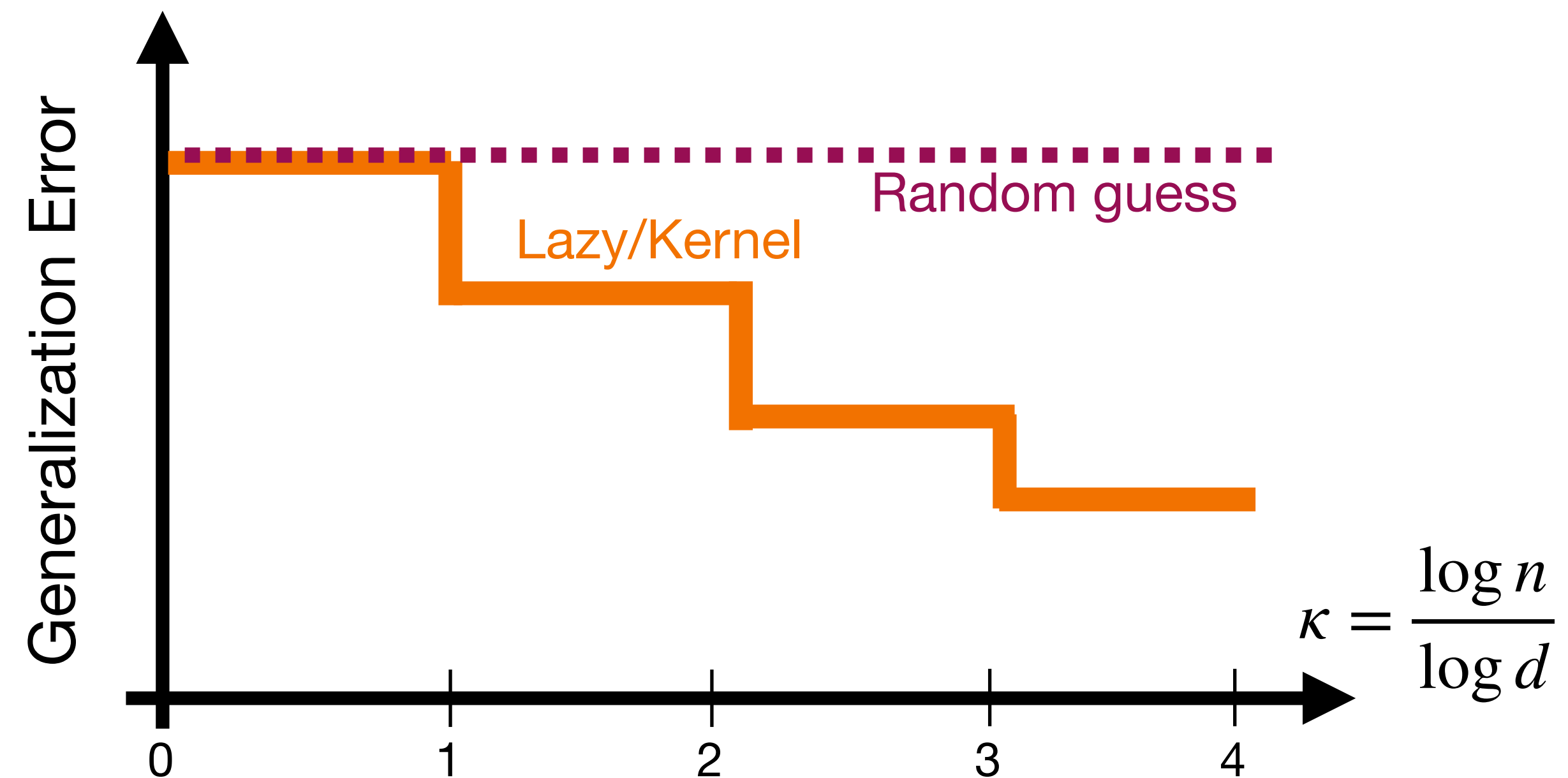




# **Learning SIGHT with Two-layer NNs**

# Lazy learning (Random Features/NTK)

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k (W^\star \mathbf{x})}{\sqrt{d^\varepsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$



# Learning with a two-layer net

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(x_{\star})}{\sqrt{d^{\varepsilon}}} \right), \mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\varepsilon}}$$

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$



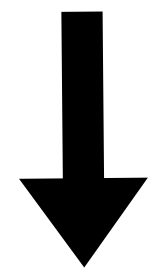
# Learning with a two-layer net

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\epsilon}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(x_{\star})}{\sqrt{d^{\epsilon}}} \right), \mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\epsilon}}$$

parameter count of  $W^{\star}$

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$



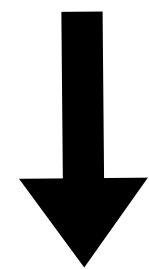
With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^{\star} + Z_2$

# Learning with a two-layer net

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \quad f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star = W^\star \mathbf{x} \in \mathbb{R}^{d^\epsilon}$$

parameter count of  $W^\star$

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$



With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

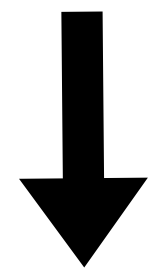
$$\hat{y} \approx \hat{W}_2 \sigma((Z_1 W^\star + Z_2) \mathbf{x}) \approx \hat{W}_2 \sigma(Z_1 x_\star + Z_3)$$

# Learning with a two-layer net

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \quad f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star = W^\star \mathbf{x} \in \mathbb{R}^{d^\epsilon}$$

parameter count of  $W^\star$

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$



With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

$$\hat{y} \approx \hat{W}_2 \sigma((Z_1 W^\star + Z_2) \mathbf{x}) \approx \hat{W}_2 \sigma(Z_1 x_\star + Z_3)$$

random feature/lazy learning but in reduced dimension  $\mathbb{R}^{d^\epsilon}$



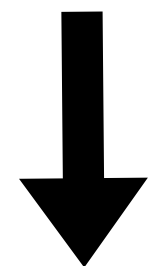
# Learning with a two-layer net

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star = W^\star \mathbf{x} \in \mathbb{R}^{d^\epsilon}$$

**parameter count of  $W^\star$**

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$



**With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$**

$$\hat{y} \approx \hat{W}_2 \sigma((Z_1 W^\star + Z_2) \mathbf{x}) \approx \hat{W}_2 \sigma(Z_1 x_\star + Z_3)$$

**random feature/lazy learning but in reduced dimension  $\mathbb{R}^{d^\epsilon}$**

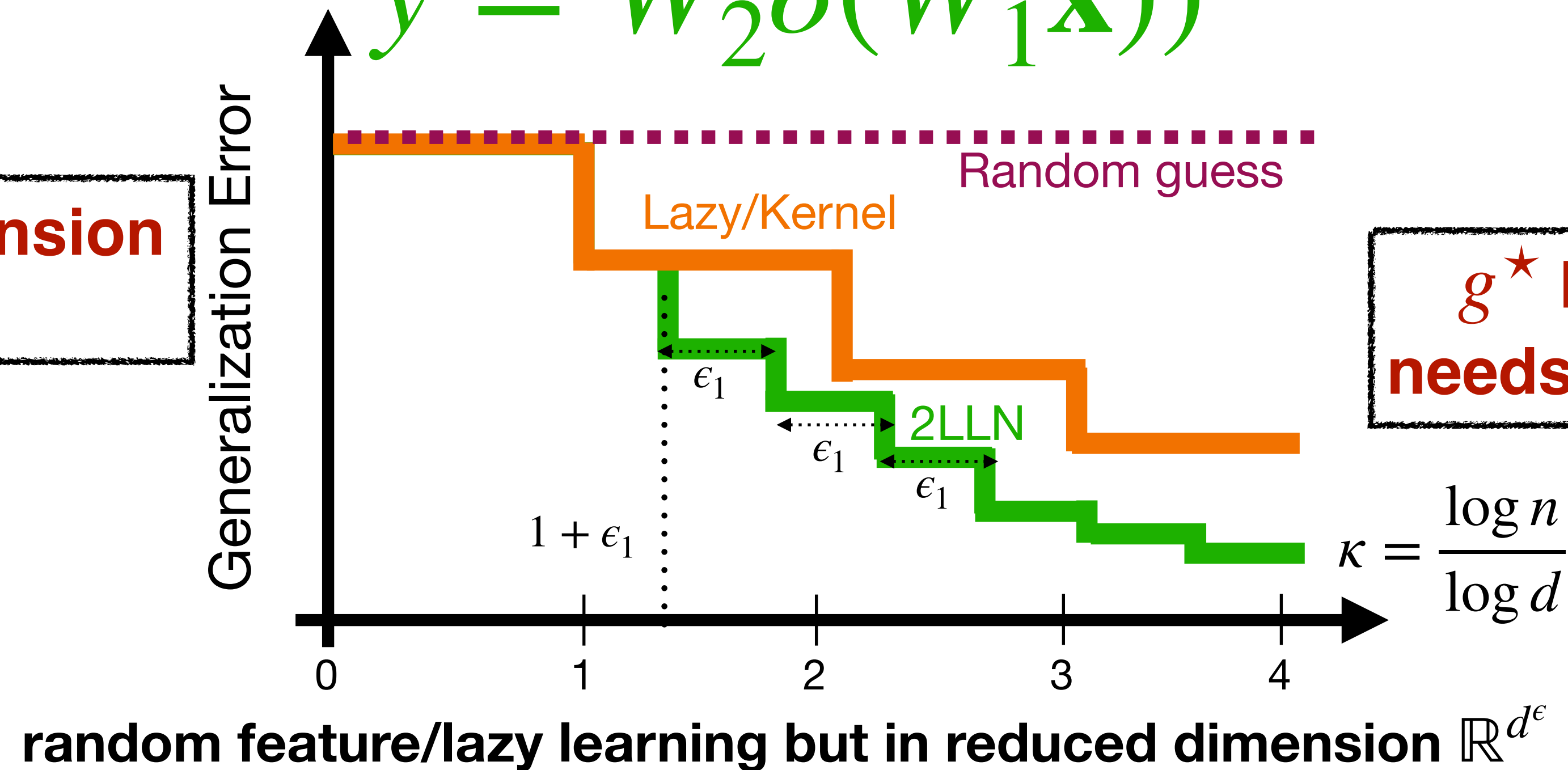
# Adaptive Learning with a two-layer net

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^{\epsilon_1}}} \right), \mathbf{x}_\star \in \mathbb{R}^{d^{\epsilon_1}}$$

$$\hat{y} = \hat{W}_2 \sigma(\hat{W}_1 \mathbf{x})$$

Only 1 level of dimension reduction



$g^\star$  has degree  $m \implies$  needs sample, width  $d^{kme_1}$





# **Learning SIGHT with Three-layer NNs**



# Idealized scenario for layer-wise training in 3-layer NNS

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star^\top} P_k (W^\star \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star(h^\star)$$

$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$\mathbf{x}_\star = W^\star \mathbf{x} \in \mathbb{R}^{d^\epsilon}$$

$\mathbf{x} \in \mathbb{R}^d$

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma (W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$

$$\mathbf{w}_3^T \in \mathbb{R}^{p_2}$$

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(\mathbf{W}_1 \mathbf{x}) \in \mathbb{R}^{p_2}$$

$$W_1 \mathbf{x} \in \mathbb{R}^{p_1}$$

$\mathbf{x} \in \mathbb{R}^d$

# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$



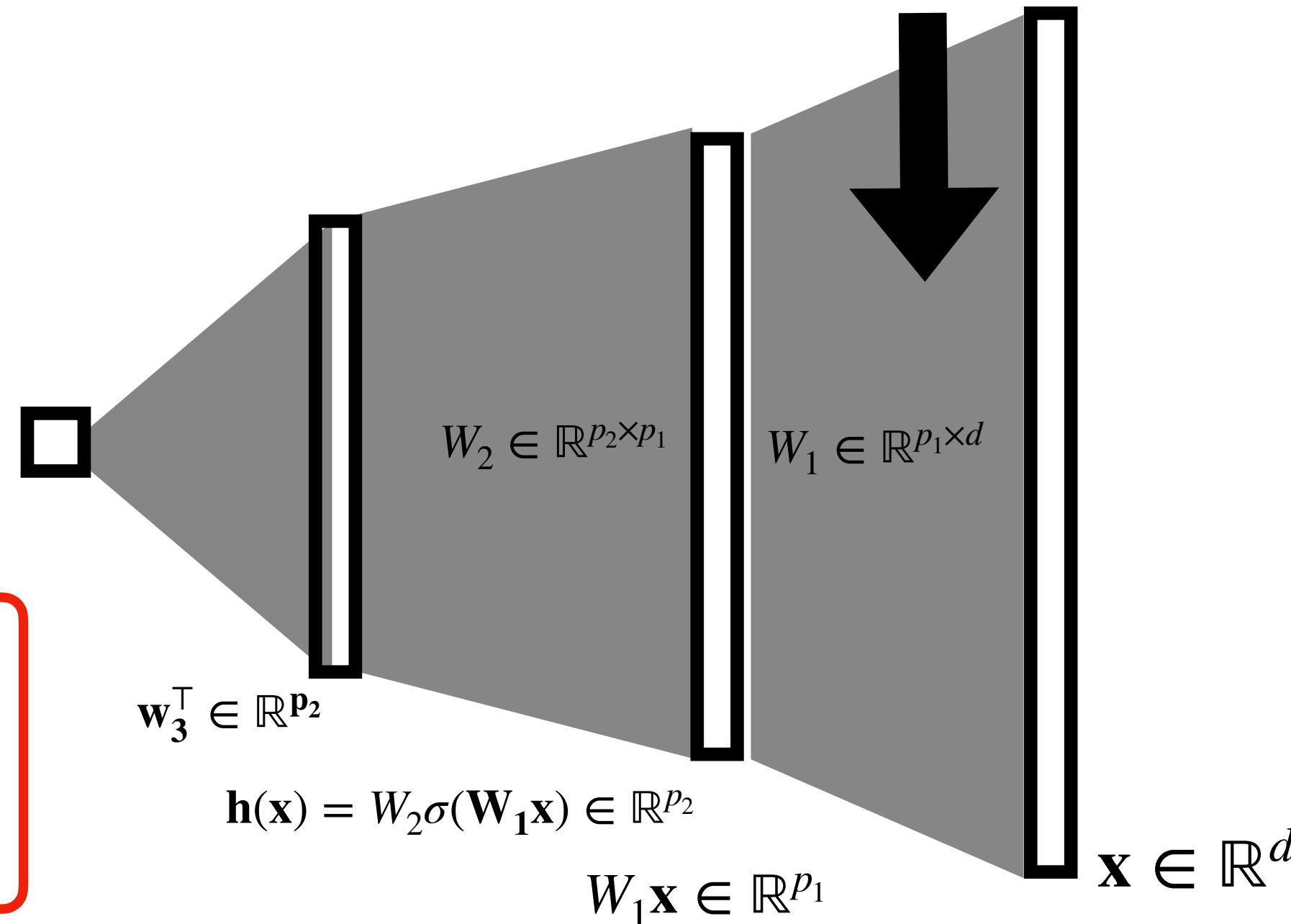
$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$$\mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\varepsilon}}$$

$$W^{\star} \in \mathbb{R}^{d^{\varepsilon} \times d}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$



# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$

$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$$\mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\varepsilon}}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$W_2 \in \mathbb{R}^{p_2 \times p_1}$$

$$W_1 \in \mathbb{R}^{p_1 \times d}$$

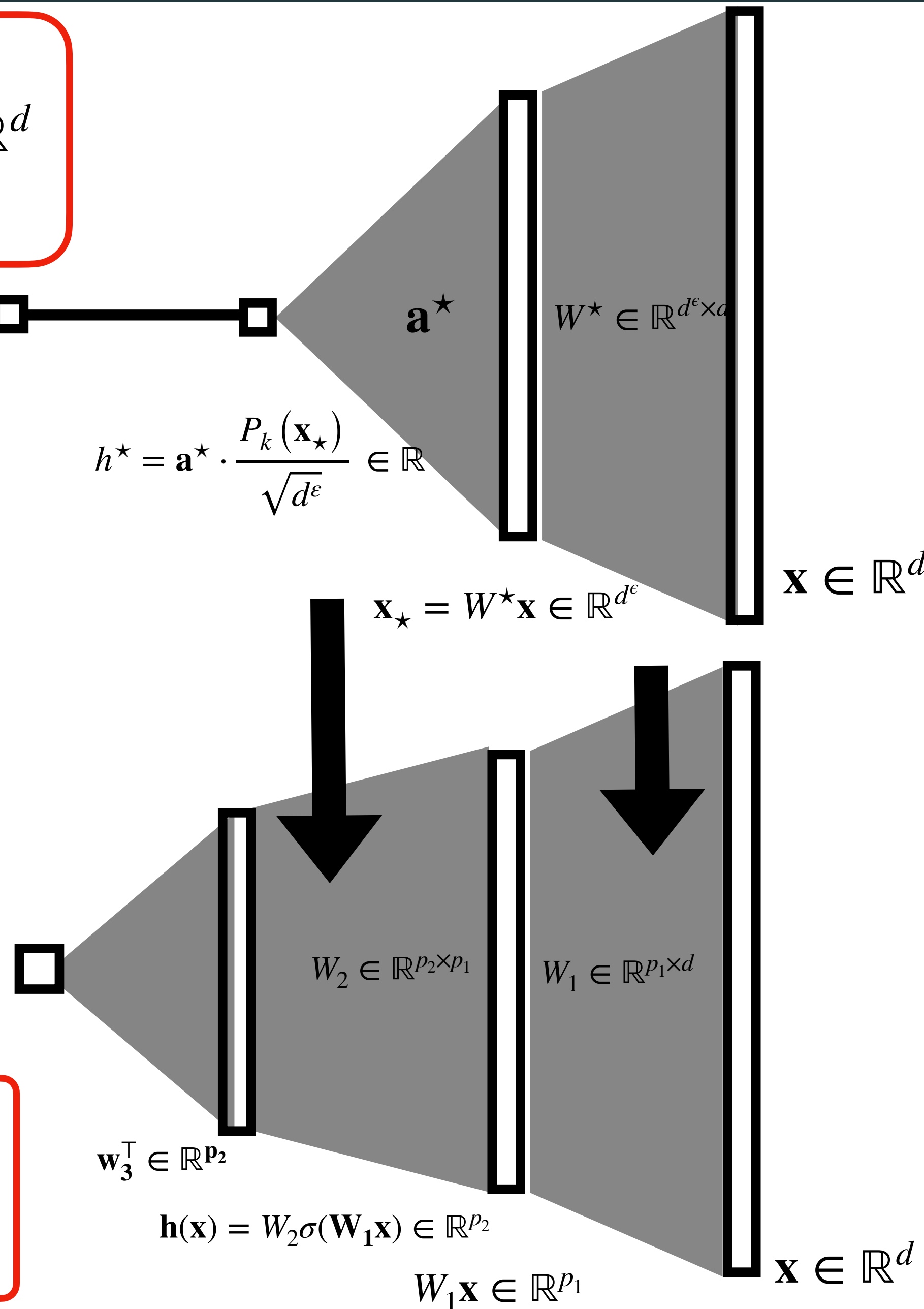
$$\mathbf{w}_3^{\top} \in \mathbb{R}^{p_2}$$

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}$$

$$W_1 \mathbf{x} \in \mathbb{R}^{p_1}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$



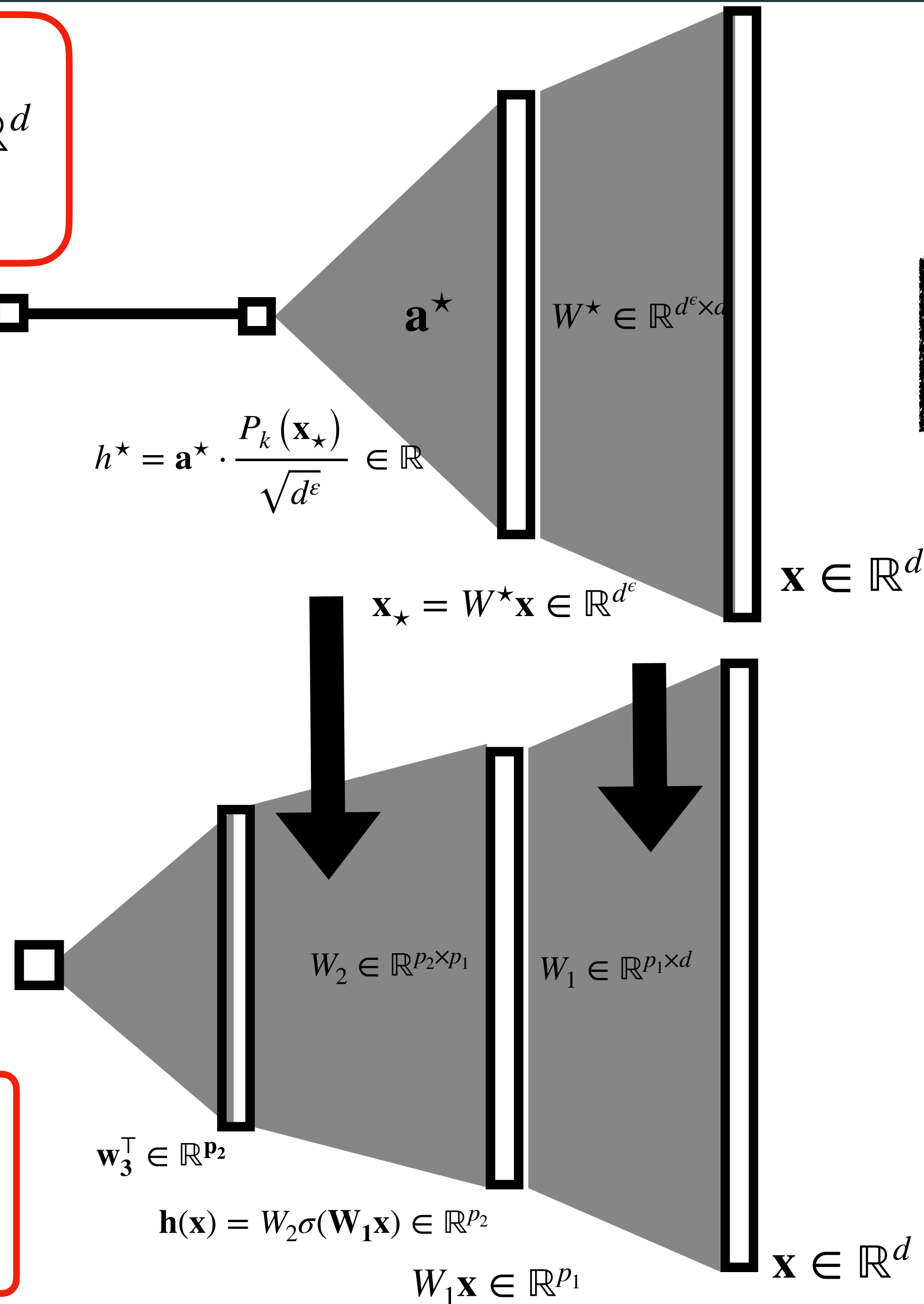
# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$

$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$W_2$  learns features not weights  $a^{\star}$



$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$



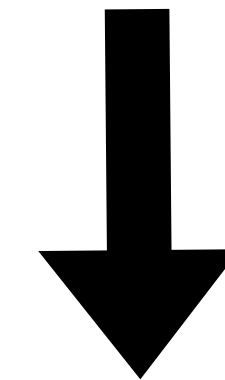
# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

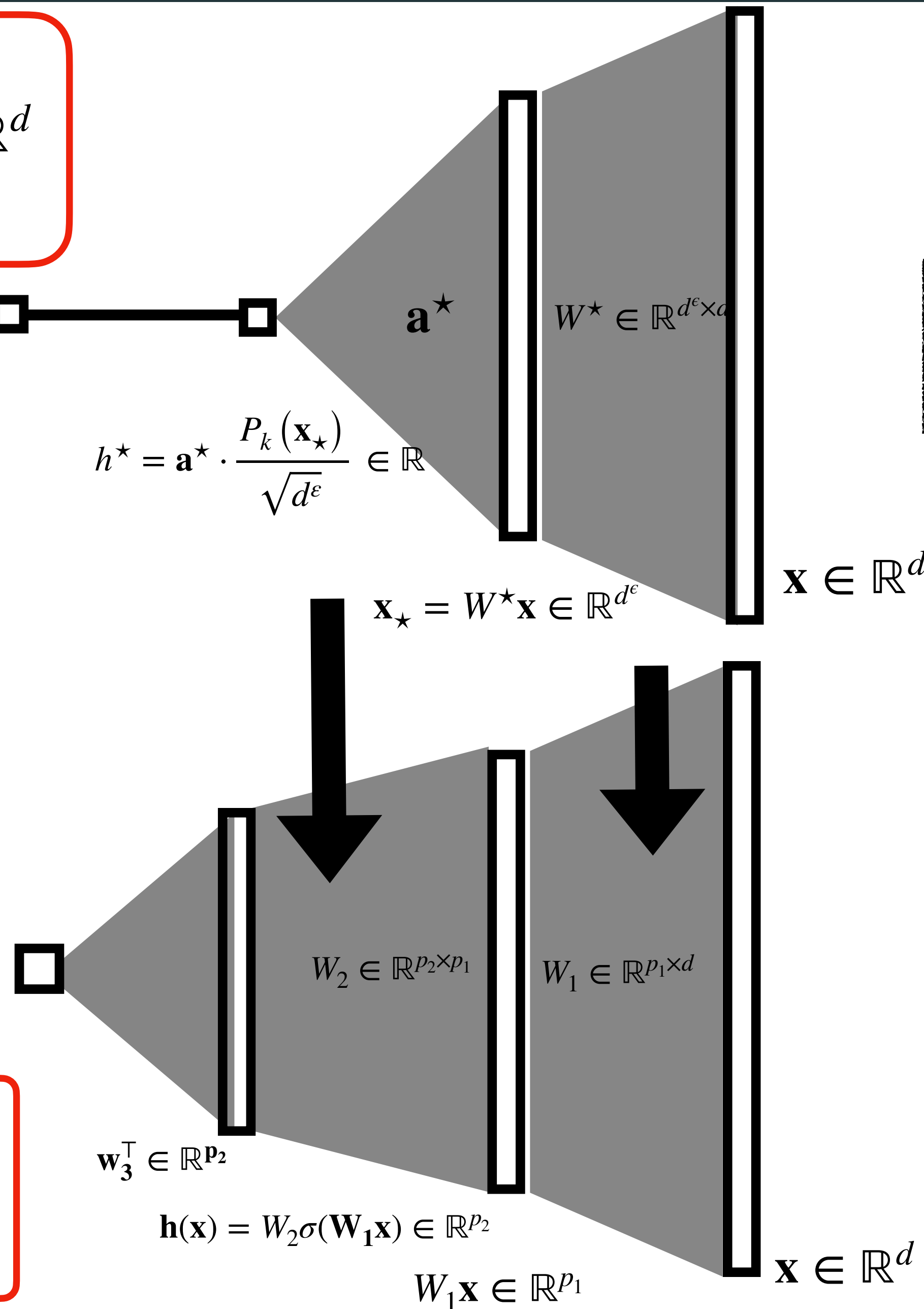
$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$

$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$W_2$  learns features not weights  $a^{\star}$



$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$



# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$

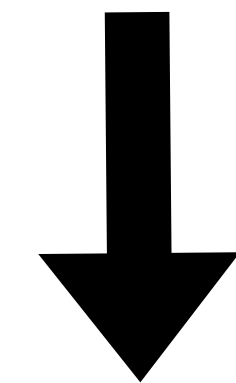
$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$$\mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\varepsilon}}$$

$$W^{\star} \in \mathbb{R}^{d^{\varepsilon} \times d}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$W_2$  learns features not weights  $a^{\star}$



$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \propto h^{\star}(\mathbf{x})$$

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$

$$\mathbf{w}_3^{\top} \in \mathbb{R}^{p_2}$$

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}$$

$$W_1 \mathbf{x} \in \mathbb{R}^{p_1}$$

$$W_1 \in \mathbb{R}^{p_1 \times d}$$

$$W_2 \in \mathbb{R}^{p_2 \times p_1}$$

$$\mathbf{x} \in \mathbb{R}^d$$

# Idealized scenario for layer-wise training in 3-layer NNS

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k(W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star}(h^{\star})$$

$$h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

$$\mathbf{x}_{\star} = W^{\star} \mathbf{x} \in \mathbb{R}^{d^{\varepsilon}}$$

$W_2$  learns features not weights  $a^{\star}$

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \propto h^{\star}(\mathbf{x})$$

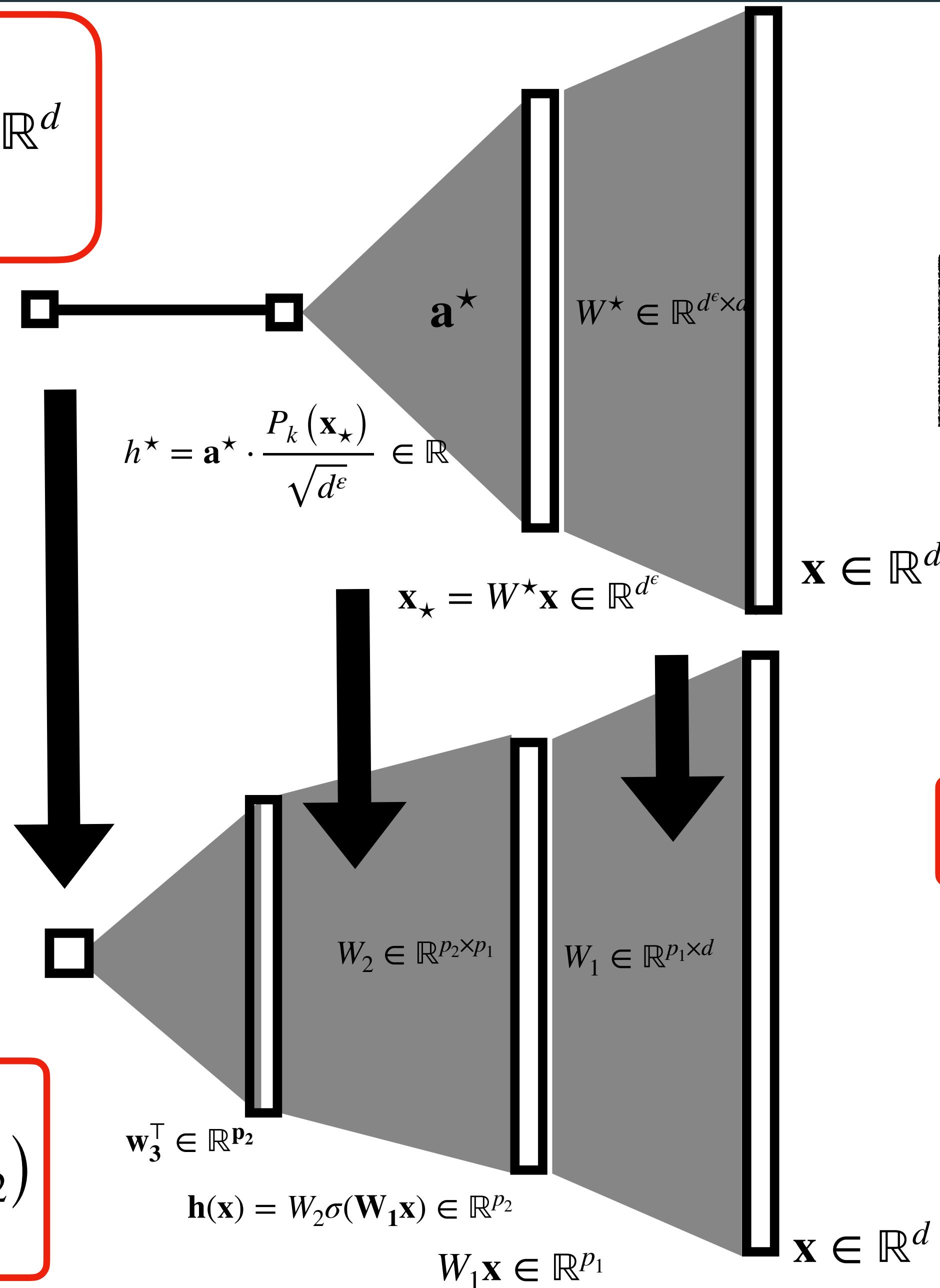
$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^{\top} \sigma(W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2)$$

$$\mathbf{w}_3^{\top} \in \mathbb{R}^{p_2}$$

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}$$

$$W_1 \mathbf{x} \in \mathbb{R}^{p_1}$$

$$\mathbf{x} \in \mathbb{R}^d$$



# Expected behavior under idealized scenario

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star\top} P_k(W^{\star}\mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star\top} P_k(x_{\star})}{\sqrt{d^{\varepsilon}}} \right), \mathbf{x}_{\star} \in \mathbb{R}^{d^{\varepsilon}}$$

$$\hat{y} = \hat{W}^3 \sigma(\hat{W}_2 \sigma(\hat{W}_1 \mathbf{x}))$$



# Expected behavior under idealized scenario

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\varepsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d^\varepsilon}$$

$$\hat{y} = \hat{\mathbf{W}}^3 \sigma(\hat{W}_2 \sigma(\hat{W}_1 \mathbf{x}))$$

With  $n = \mathcal{O}(d^{1+\varepsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

$$\hat{y} \approx \hat{\mathbf{w}}_3^\top \sigma(\hat{W}_2 \sigma(Z_1 x_\star + Z_2))$$

# Expected behavior under idealized scenario

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^{\epsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d^\epsilon}$$

$$\hat{y} = \hat{\mathbf{W}}^3 \sigma(\hat{W}_2 \sigma(\hat{W}_1 \mathbf{x}))$$

With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

With each  $n = \mathcal{O}(d^\epsilon)$  data,  $\hat{W}_2 \sigma(\cdot)$  learns a higher-degree component of  $P_k$

$$\hat{y} \approx \hat{\mathbf{w}}_3^\top \sigma(\hat{W}_2 \sigma(Z_1 x_\star + Z_2))$$

# Expected behavior under idealized scenario

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^{\epsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(x_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d^\epsilon}$$

$$\hat{y} = \hat{\mathbf{W}}^3 \sigma(\hat{W}_2 \sigma(\hat{W}_1 \mathbf{x}))$$

With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

With each  $n = \mathcal{O}(d^\epsilon)$  data,  $\hat{W}_2 \sigma(\cdot)$  learns a higher-degree component of  $P_k$

$$\hat{y} \approx \hat{\mathbf{w}}_3^\top \sigma(\hat{W}_2 \sigma(Z_1 x_\star + Z_2))$$

$$\approx \hat{\mathbf{w}}_3^\top \sigma(Z_1 h^\star + Z_2)$$

# Expected behavior under idealized scenario

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^{\varepsilon_1}}} \right), \mathbf{x} \in \mathbb{R}^d$$

$$f^\star(\mathbf{x}) = g^\star(h^\star), h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\varepsilon}} \in \mathbb{R}$$

$$\hat{y} = \hat{\mathbf{W}}^3 \sigma(\hat{W}_2 \sigma(\hat{W}_1 \mathbf{x}))$$

With  $n = \mathcal{O}(d^{1+\epsilon})$  data,  $\hat{W}_1$  can learn the features with GD :  $\hat{W}_1 \approx Z_1 W^\star + Z_2$

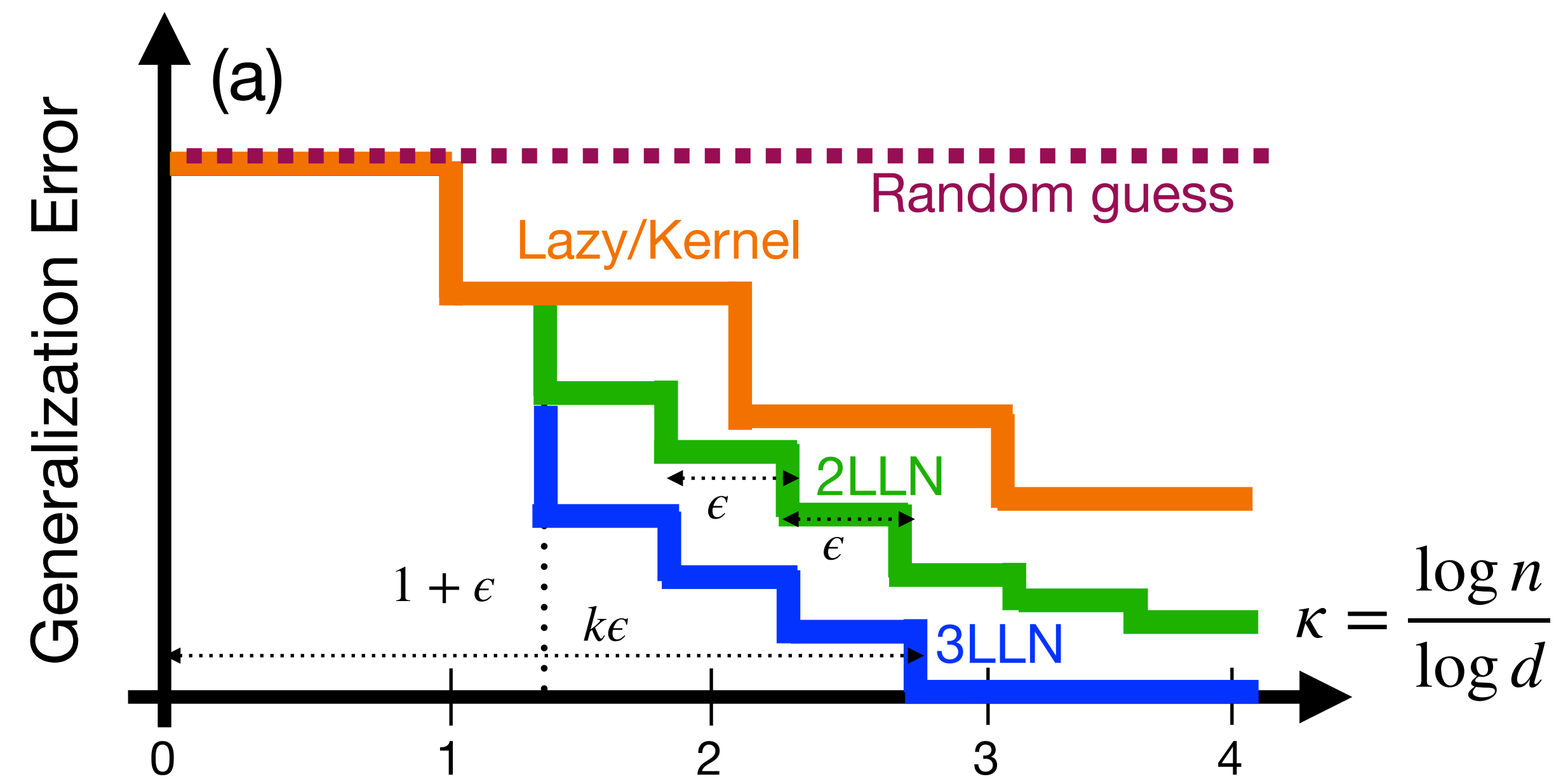
With each  $n = \mathcal{O}(d^\epsilon)$  data,  $\hat{W}_2 \sigma(\cdot)$  learns a higher-degree component of  $P_k$

$$\hat{y} \approx \hat{\mathbf{w}}_3^\top \sigma(\hat{W}_2 \sigma(Z_1 x_\star + Z_2))$$

$$\approx \hat{\mathbf{w}}_3^\top \sigma(Z_1 h^\star + Z_2)$$

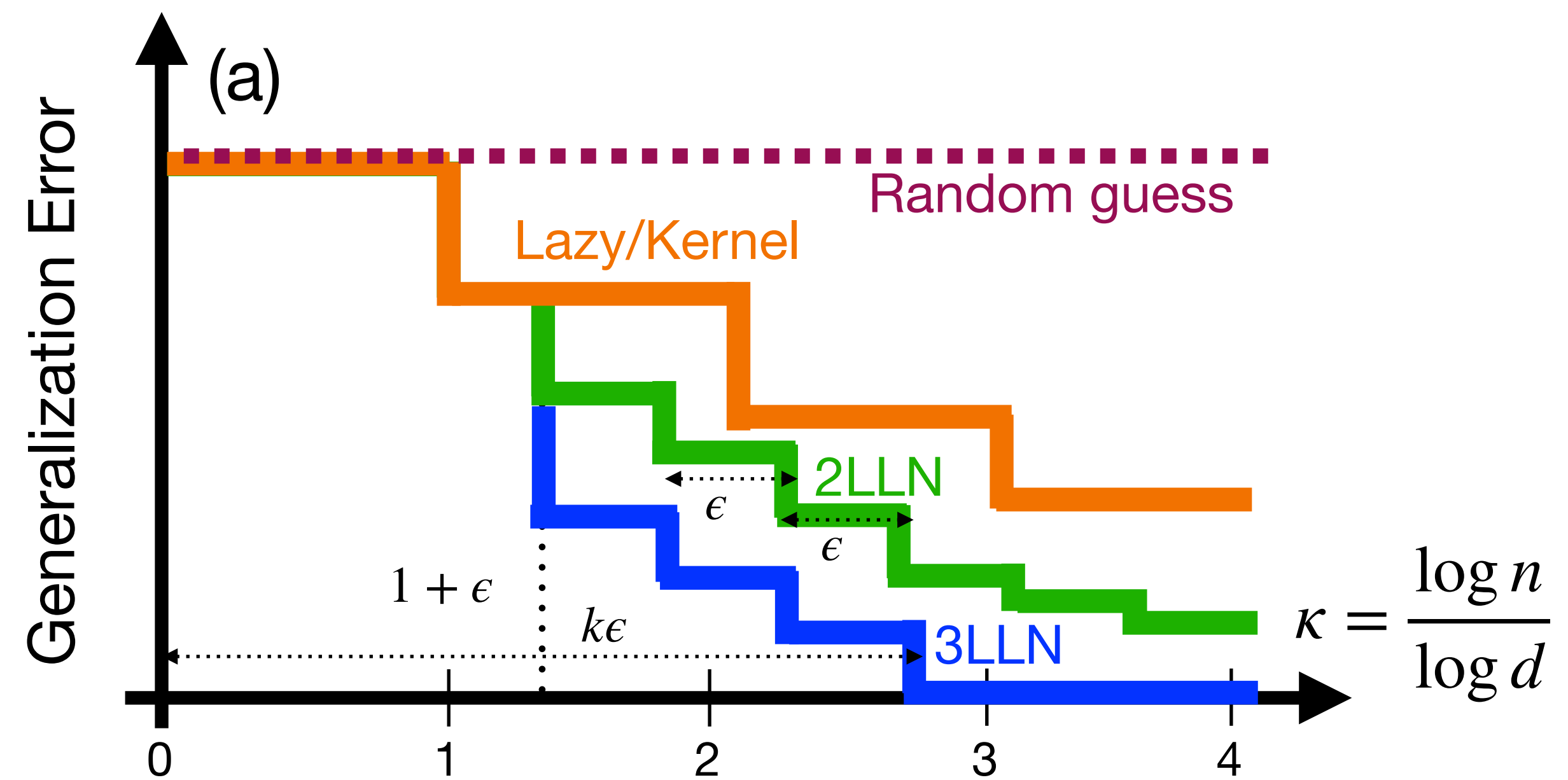


# Two levels of dimension reduction



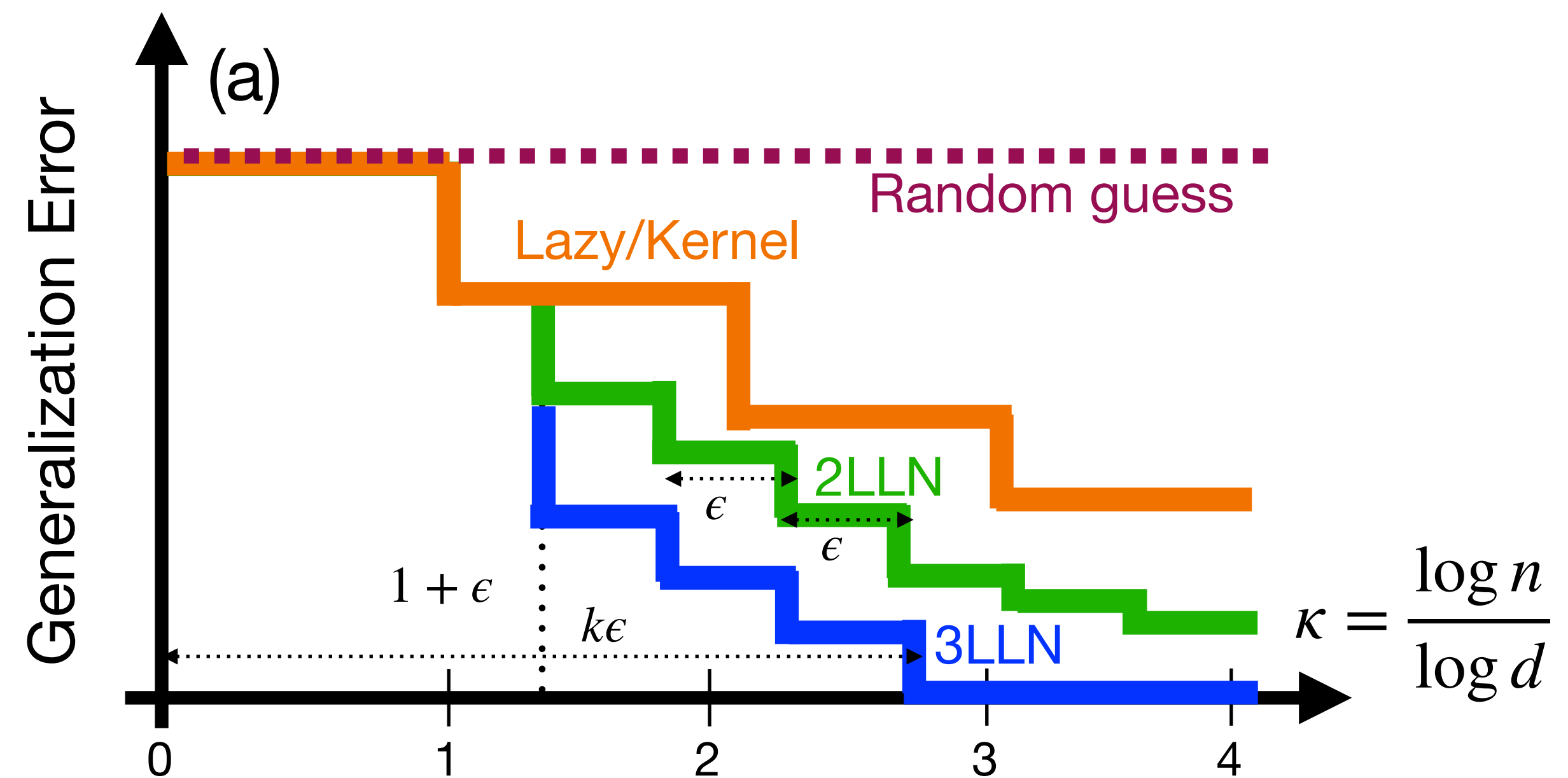
# Two levels of dimension reduction

$$f^\star(\mathbf{x}) = \mathbf{g}^\star \left( \frac{\mathbf{a}^{\star\top} \mathbf{P}_k (\mathbf{W}^\star \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$



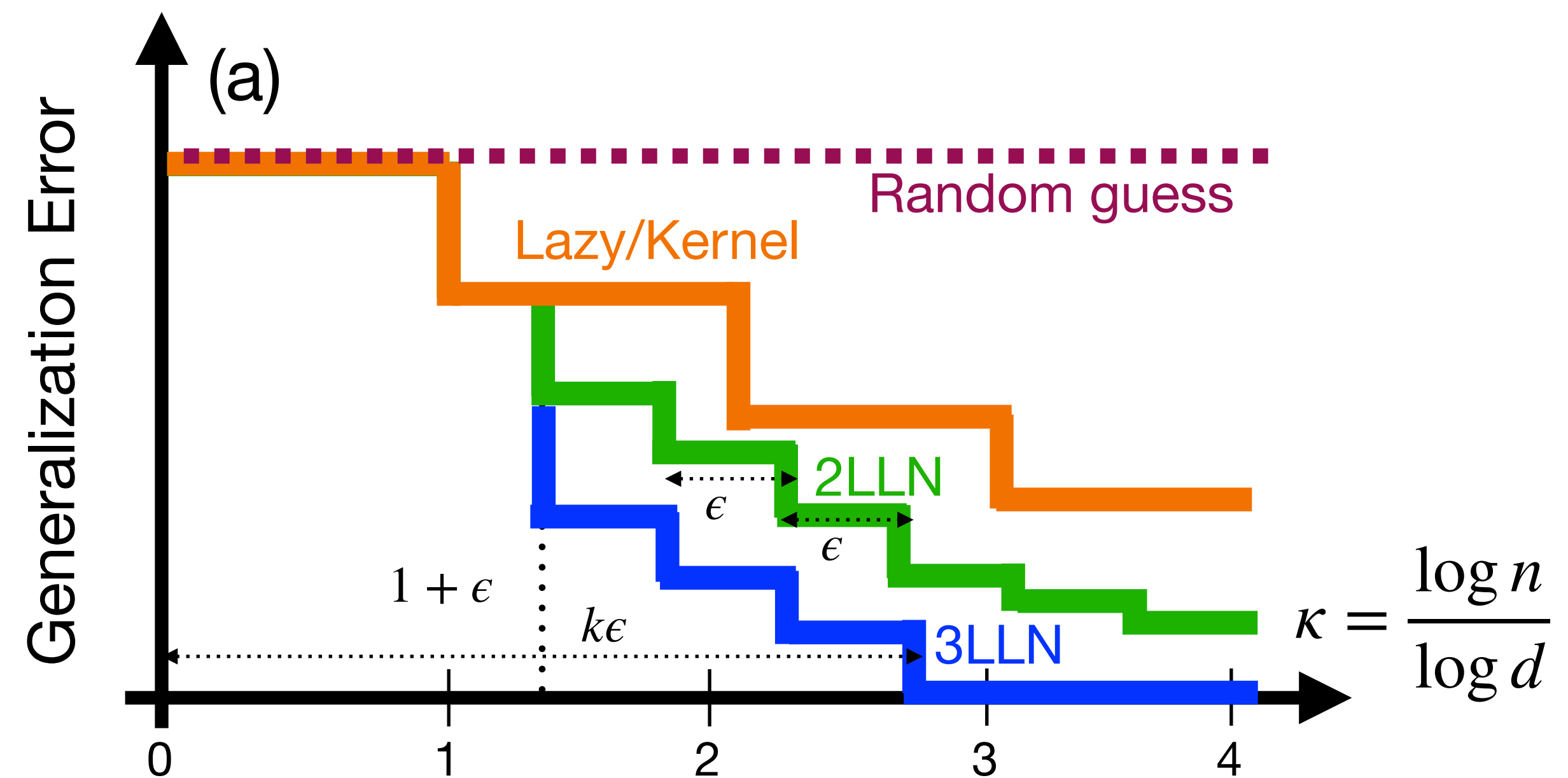
# Two levels of dimension reduction

$$f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{W}^* \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \quad \longrightarrow$$



# Two levels of dimension reduction

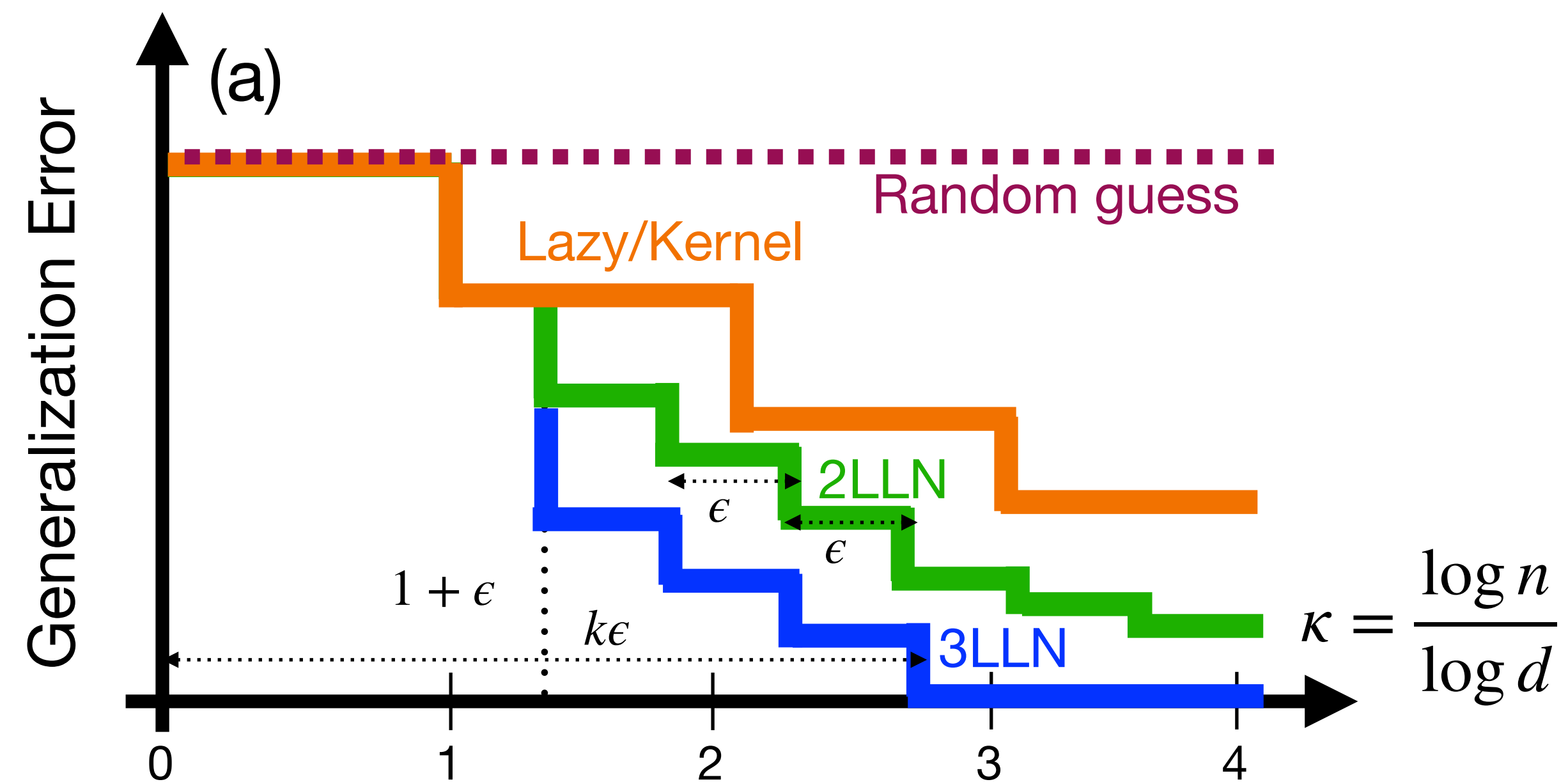
$$f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{W}^* \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \quad \longrightarrow \quad f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d_\epsilon}$$





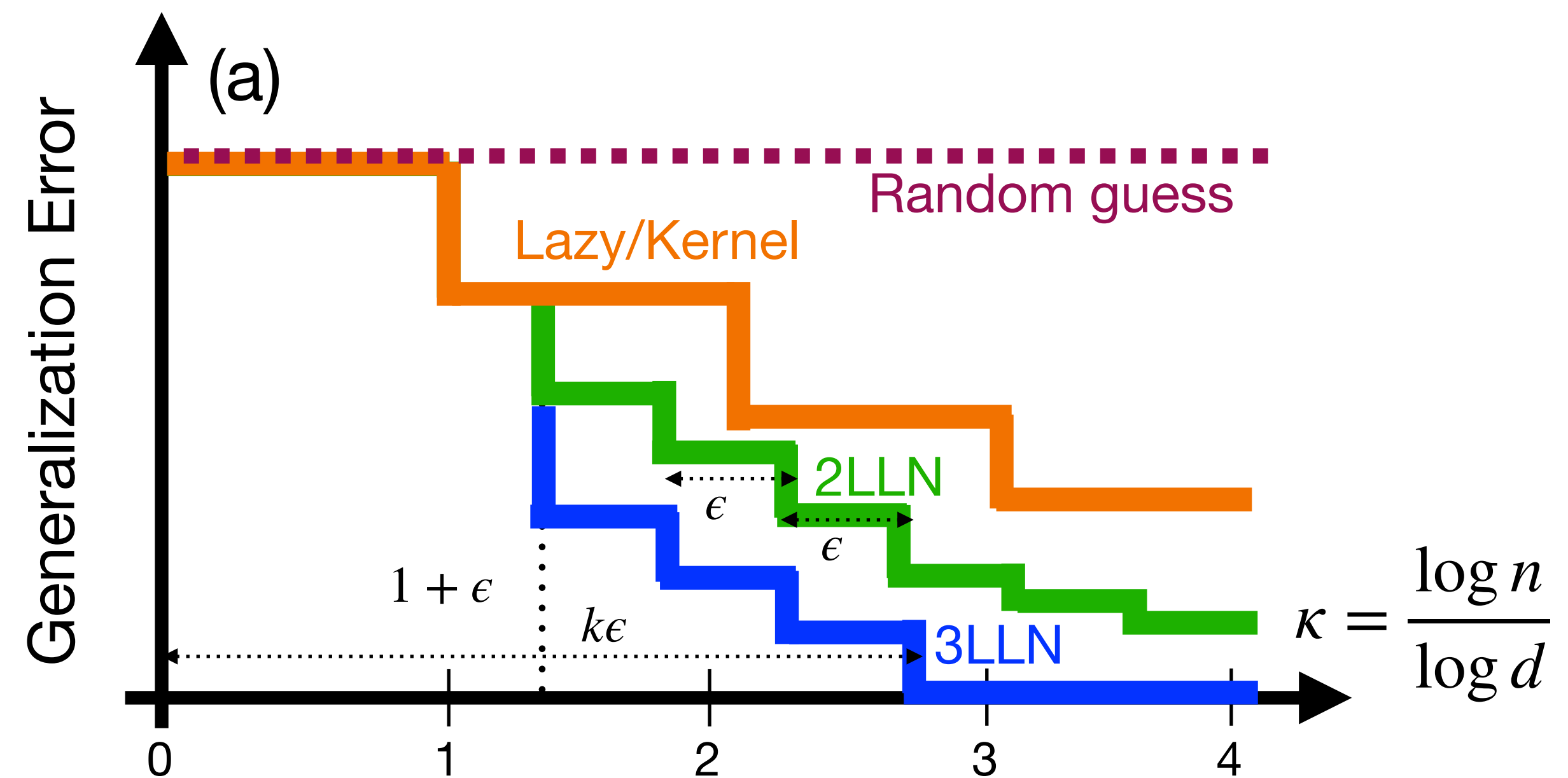
# Two levels of dimension reduction

$$f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{W}^* \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \quad \longrightarrow \quad f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d_\epsilon} \quad \longrightarrow$$

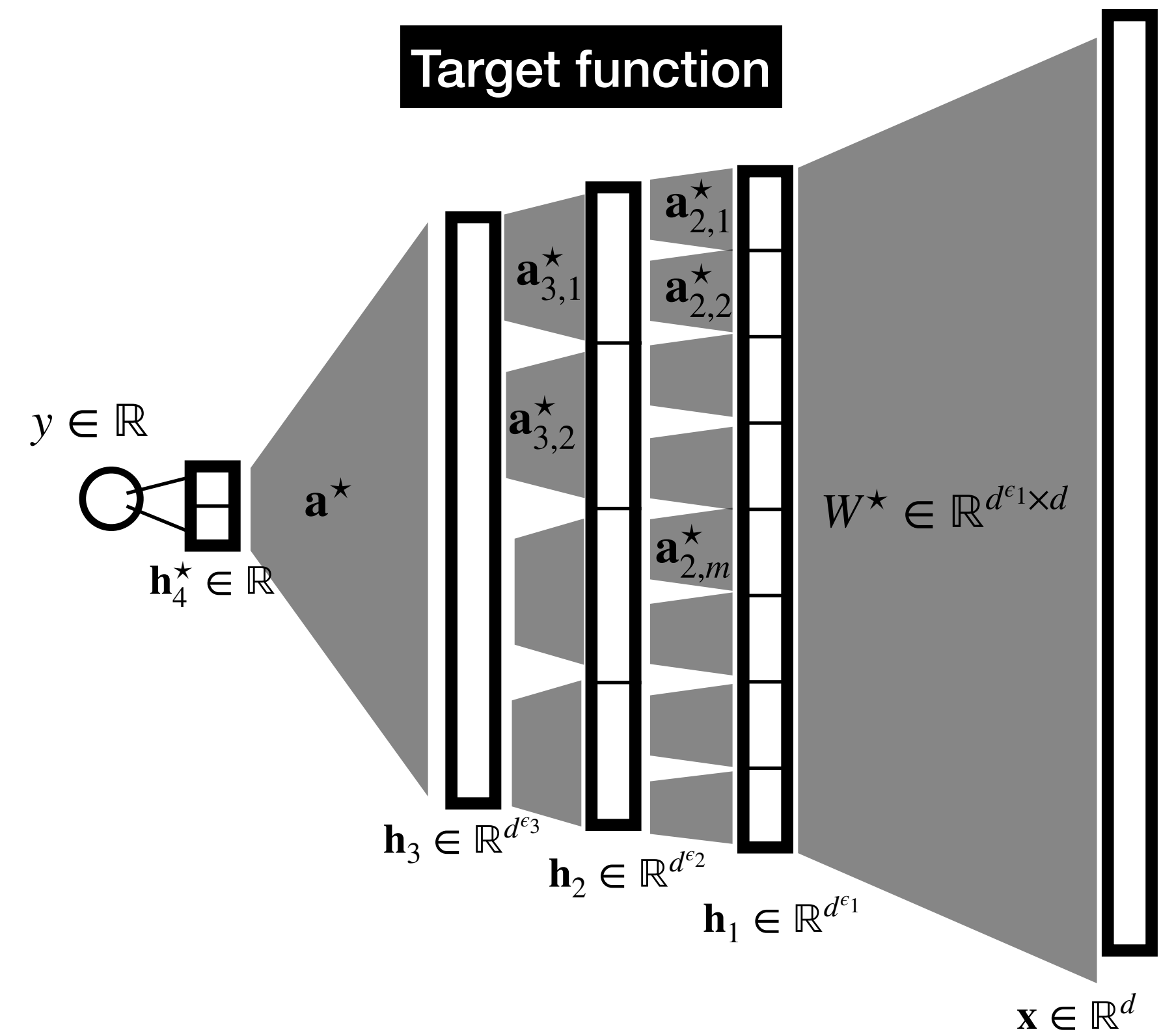
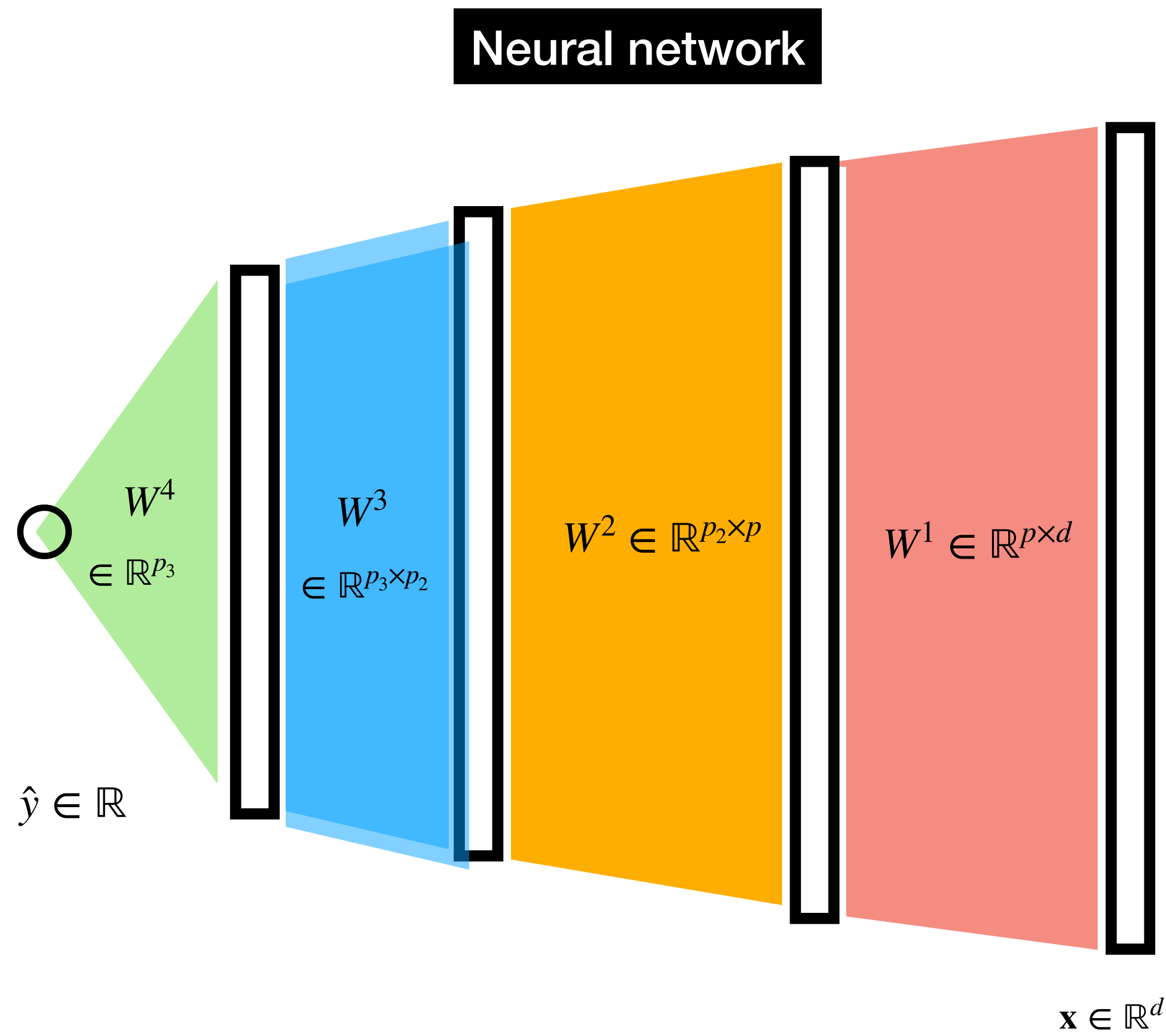


# Two levels of dimension reduction

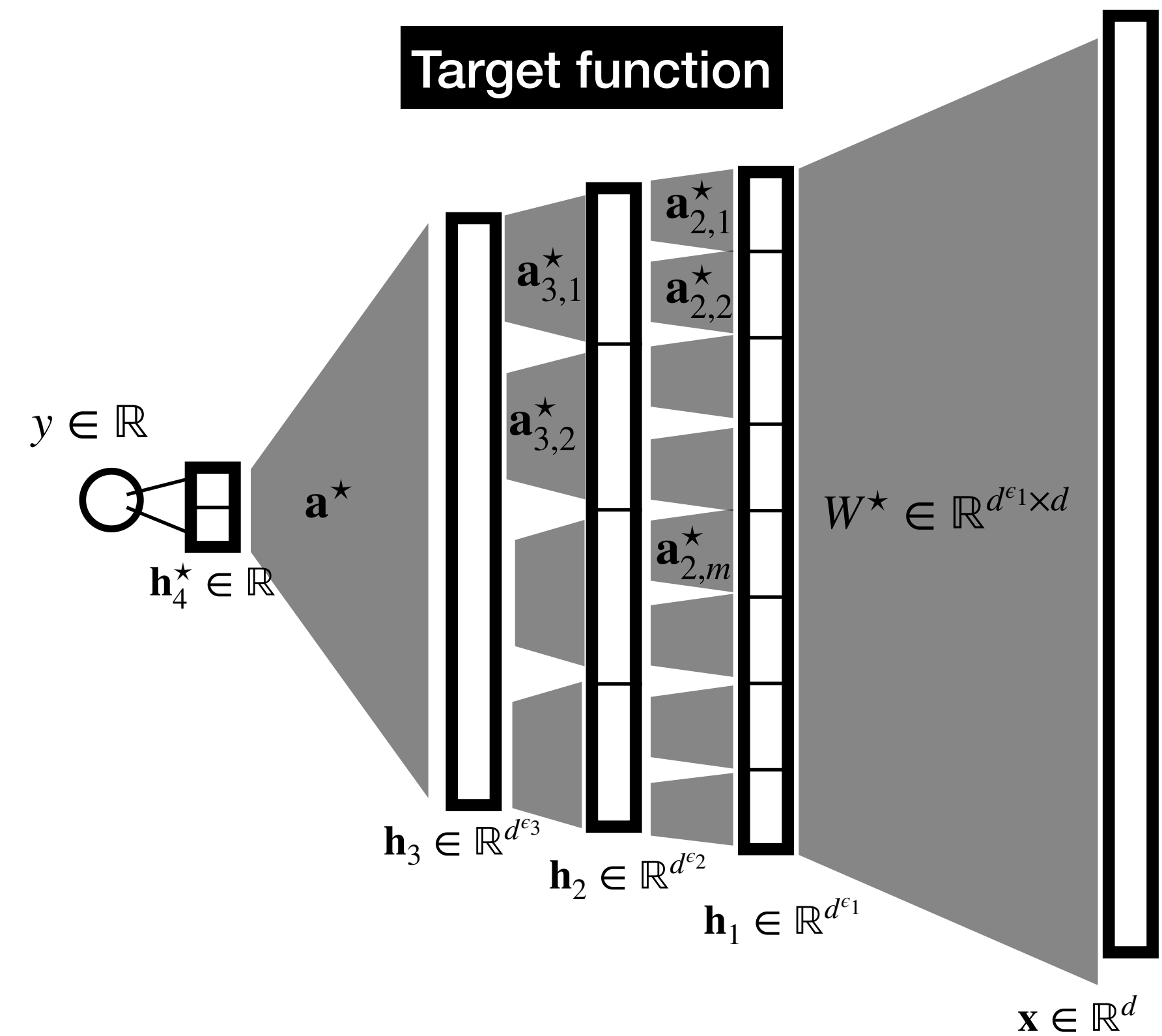
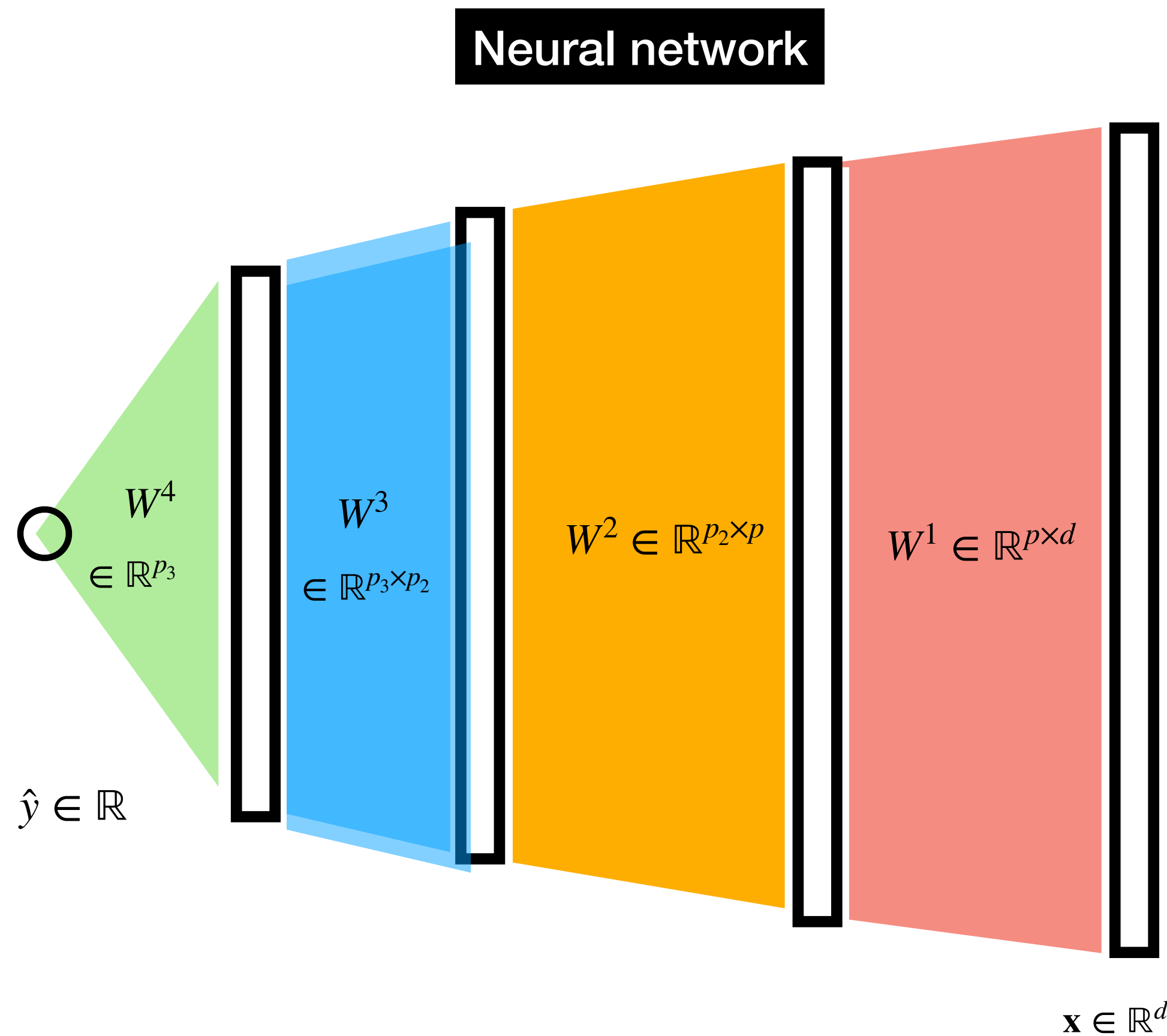
$$f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{W}^* \mathbf{x})}{\sqrt{d^\epsilon}} \right), \mathbf{x} \in \mathbb{R}^d \longrightarrow f^*(\mathbf{x}) = \mathbf{g}^* \left( \frac{\mathbf{a}^{*\top} \mathbf{P}_k (\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \right), \mathbf{x}_\star \in \mathbb{R}^{d_\epsilon} \longrightarrow \boxed{f^*(\mathbf{x}) = g^* (h^*)}$$



# Deep version



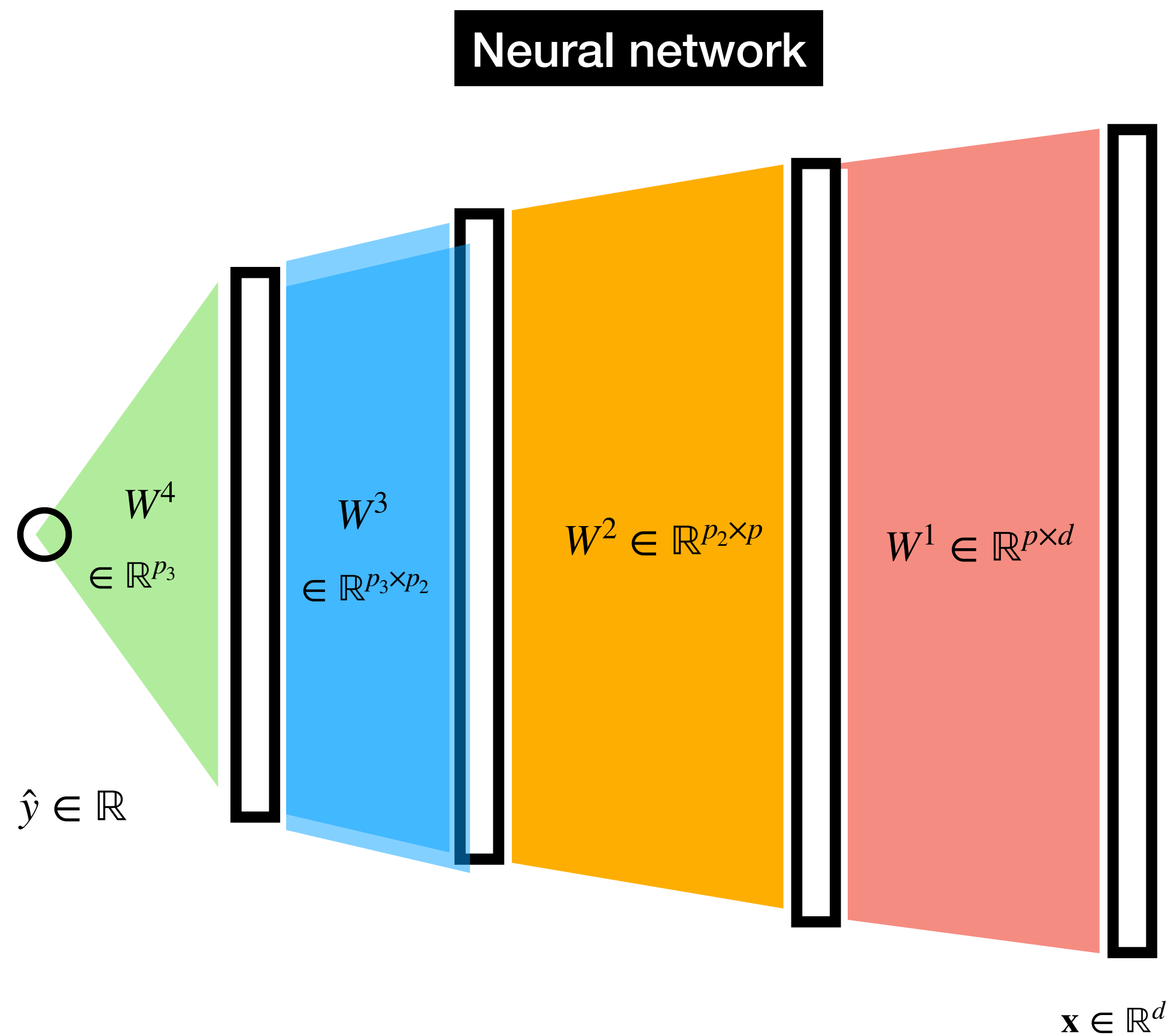
# Deep version



**Tree-structure maintains independence of features**

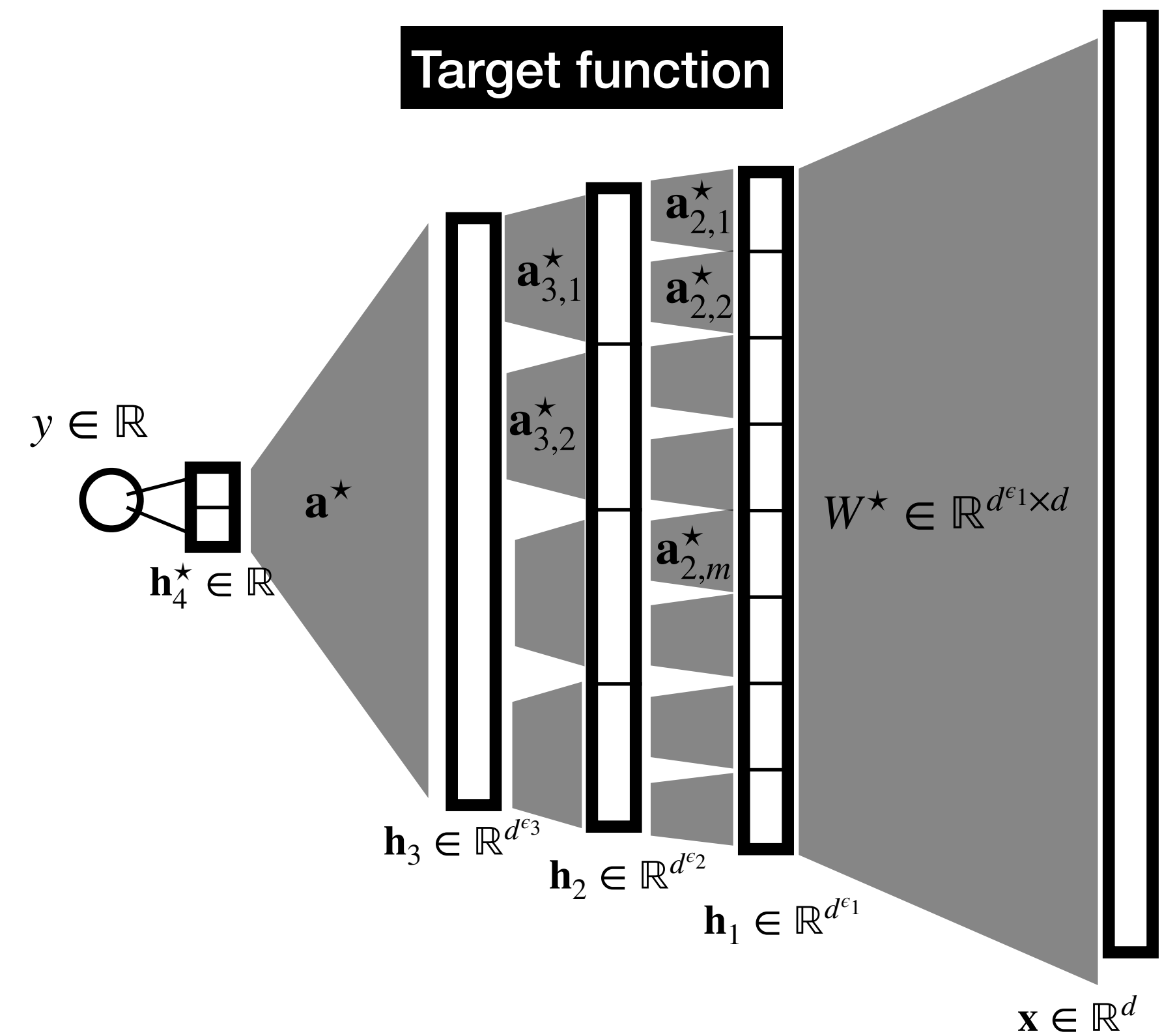


# Deep version



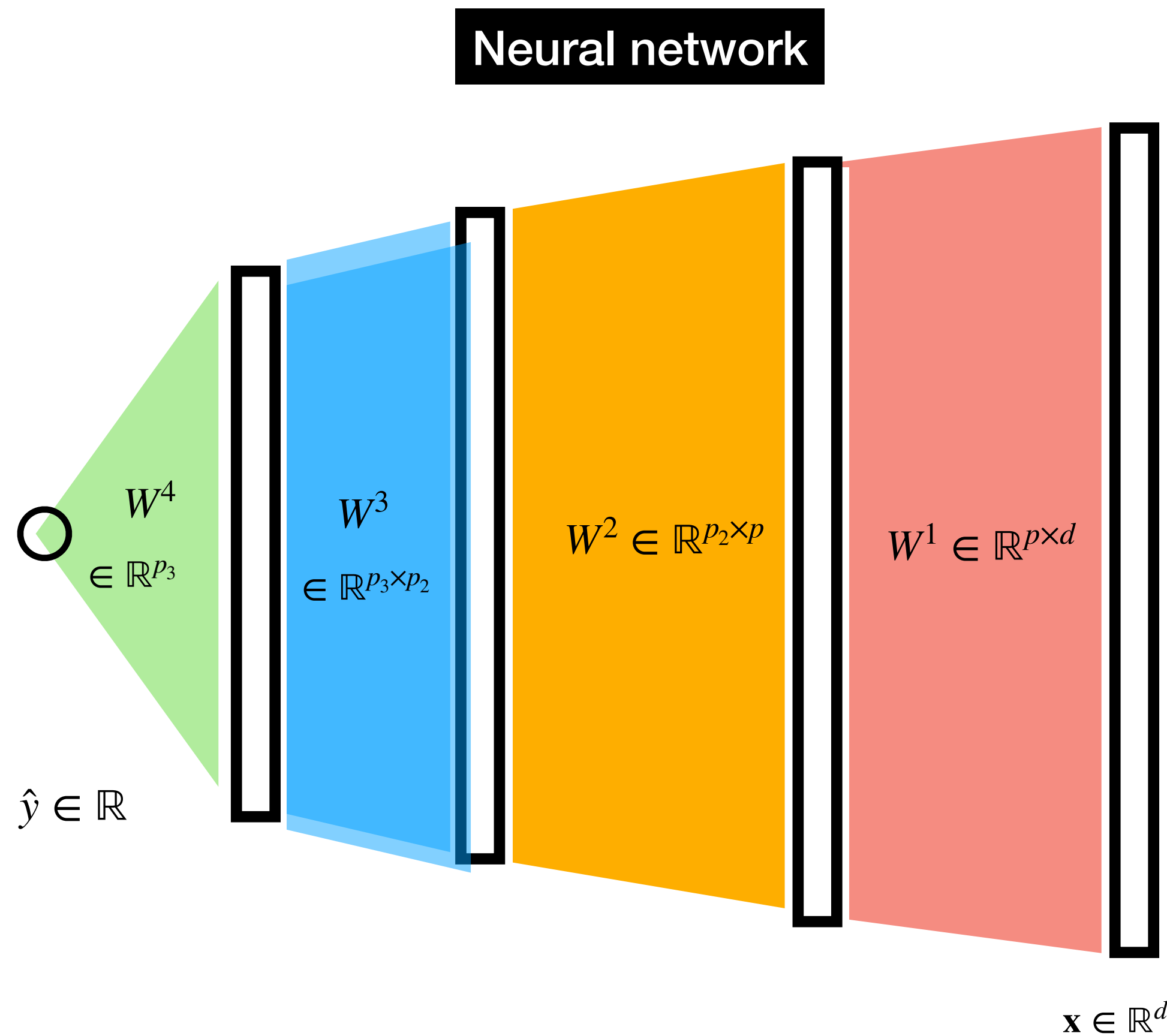
**Iterative dimensionality reduction**

$$d^{\epsilon_1} \rightarrow d^{\epsilon_2} \rightarrow d^{\epsilon_3}, \dots, \rightarrow 1$$



**Tree-structure maintains independence of features**

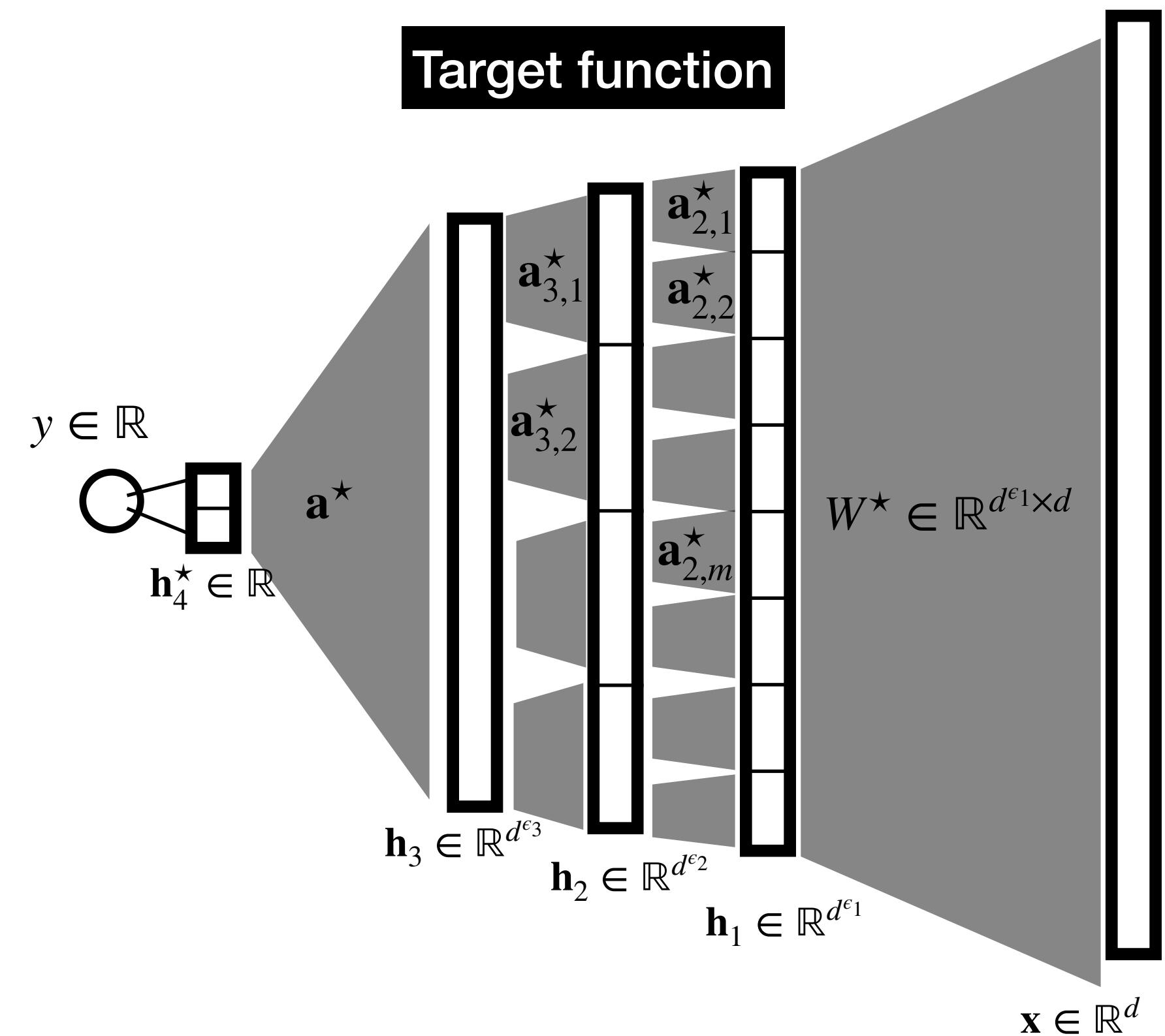
# Deep version



Iterative dimensionality reduction

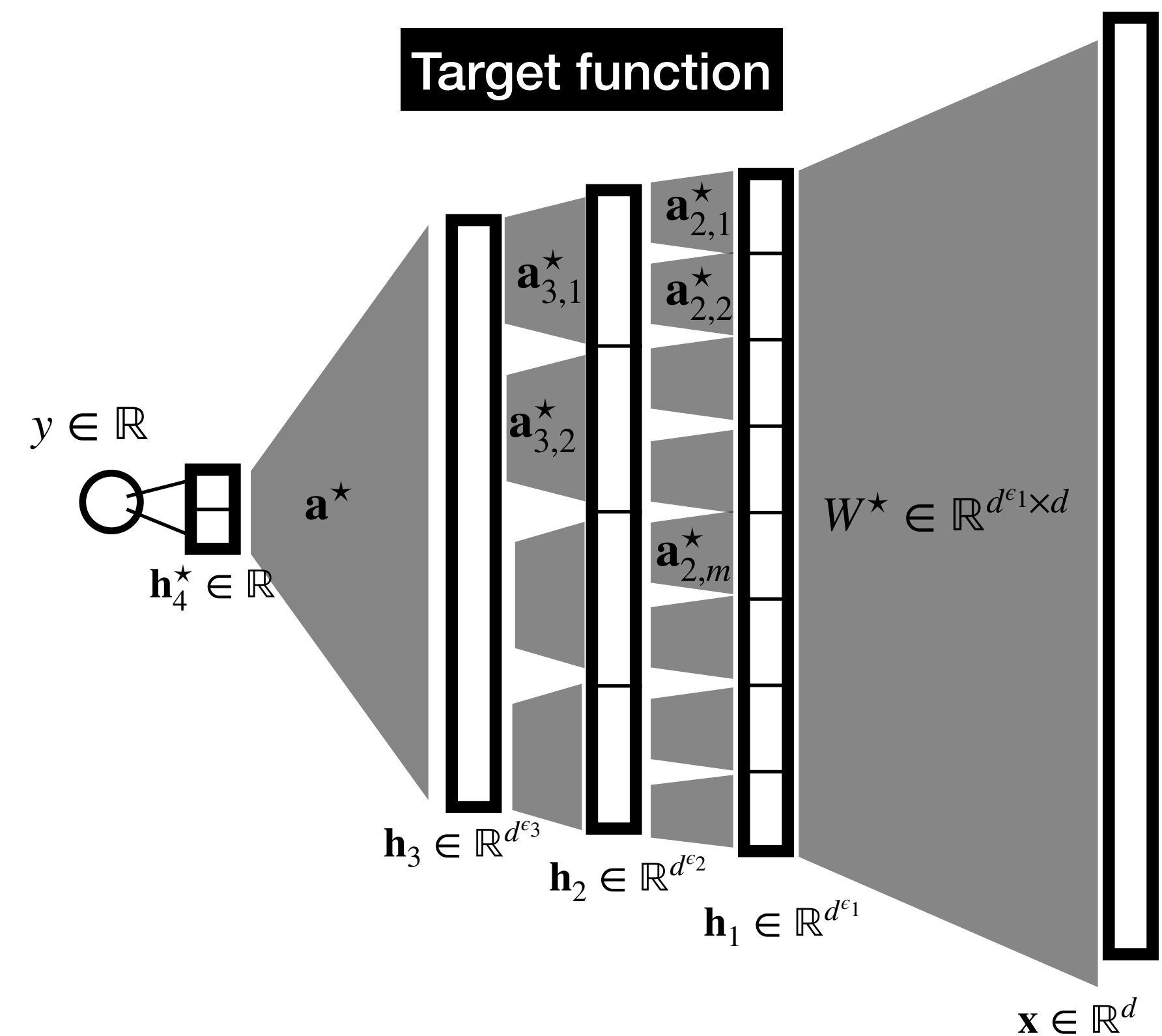
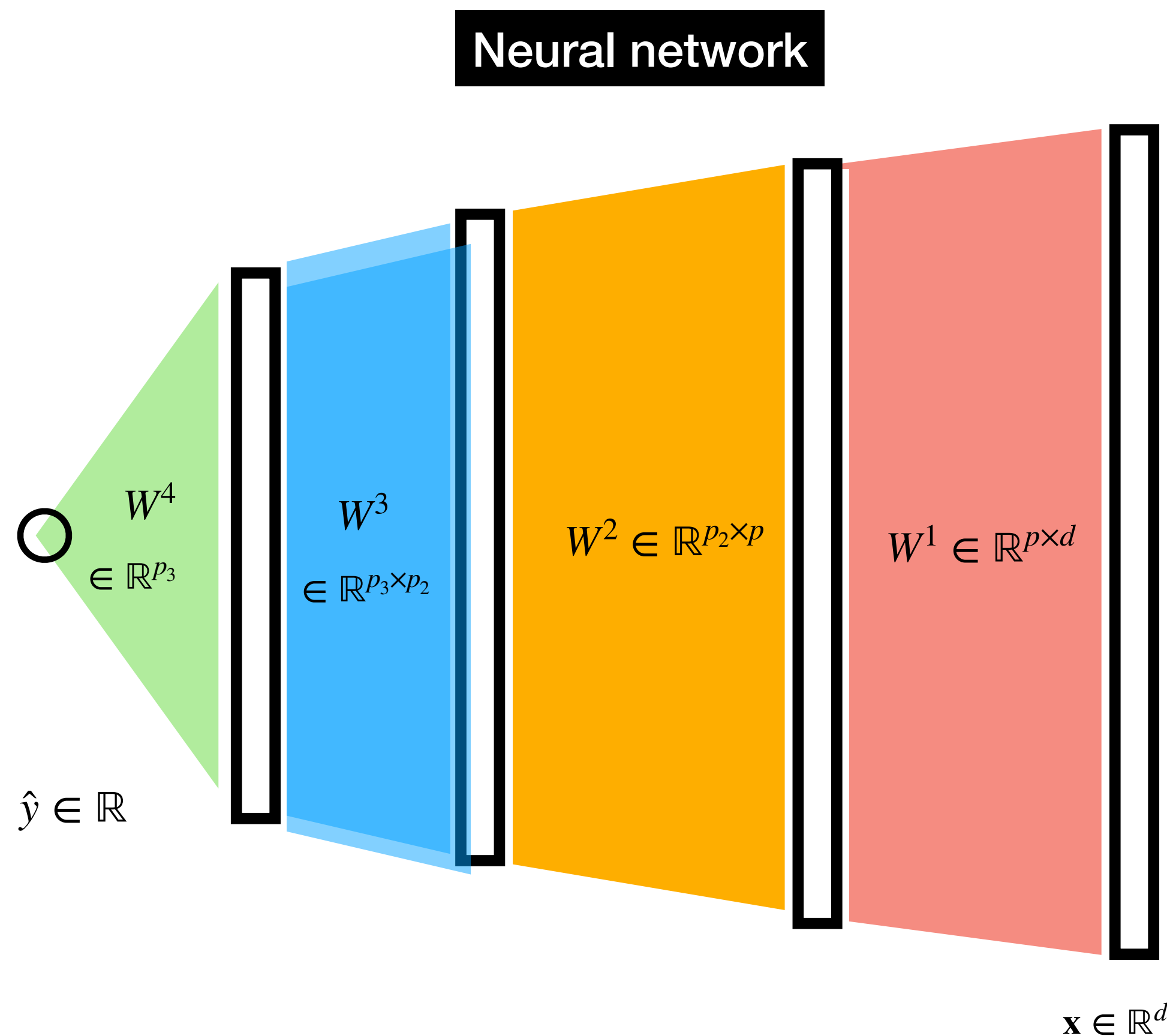
$$d^{\epsilon_1} \rightarrow d^{\epsilon_2} \rightarrow d^{\epsilon_3}, \dots, \rightarrow 1$$

**L-layer network**



Tree-structure maintains independence of features

# Deep version



**Iterative dimensionality reduction**

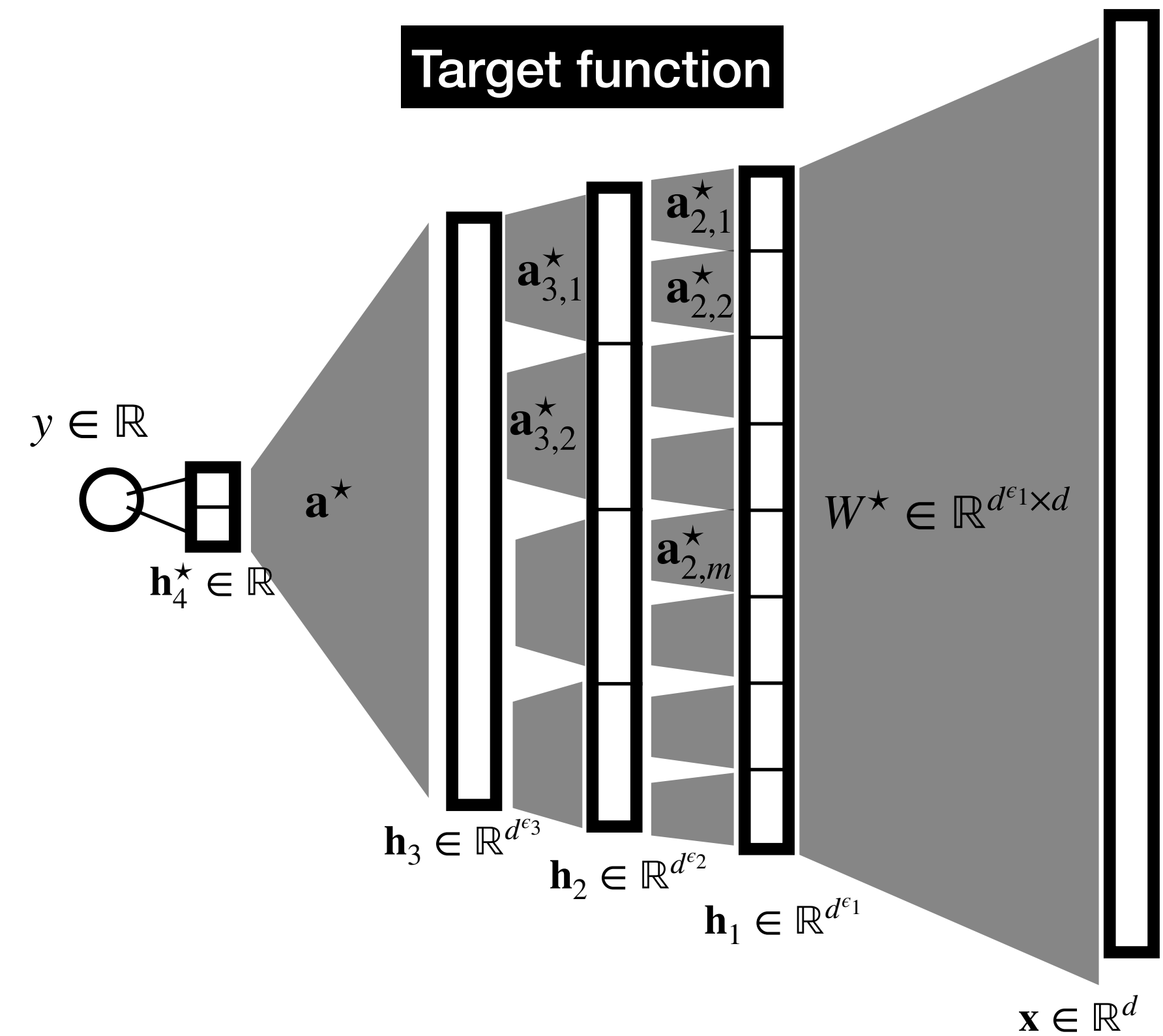
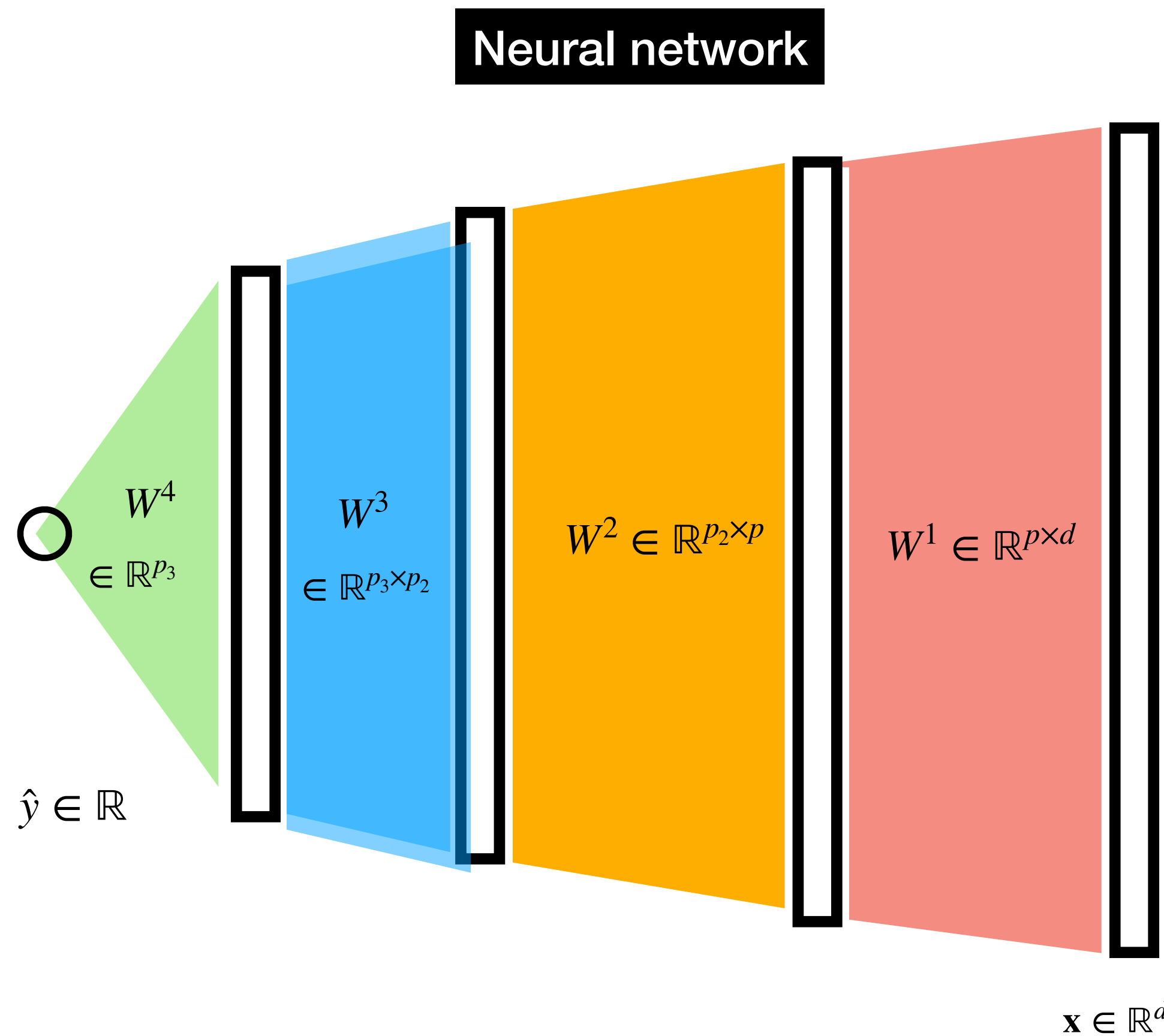
$$d^{\epsilon_1} \rightarrow d^{\epsilon_2} \rightarrow d^{\epsilon_3}, \dots, \rightarrow 1$$

**L-layer network**



**Tree-structure maintains independence of features**

# Deep version



**Iterative dimensionality reduction**

$$d^{\epsilon_1} \rightarrow d^{\epsilon_2} \rightarrow d^{\epsilon_3}, \dots, \rightarrow 1$$

**L-layer network**



**Tree-structure maintains independence of features**

**L-1 levels of dimension reduction**

# Conditions for learnability



# Conditions for learnability

## Information exponent

$$\text{IE}(\ell) = \inf\{k : \|\mathbb{E}[(\mathbf{x})^{\otimes k} f^{\star}(\mathbf{x})]\|_F \neq 0\}$$

# Conditions for learnability

**Information exponent**



$$\text{IE}(\ell) = \inf\{k : \|\mathbb{E}[(\mathbf{x})^{\otimes k} f^{\star}(\mathbf{x})]\|_F \neq 0\}$$

# Conditions for learnability

**Information exponent**



**Compositional  
Information exponent**

$$\text{IE}(\ell) = \inf\{k : \|\mathbb{E}[(\mathbf{x})^{\otimes k} f^{\star}(\mathbf{x})]\|_F \neq 0\}$$

$$\text{CIE}(\ell) = \inf\{k : \|\mathbb{E}[(h_{\ell}(\mathbf{x}))^{\otimes k} f^{\star}(\mathbf{x})]\|_F = \Theta(1)\}$$

# Conditions for learnability

**Information exponent**



**Compositional  
Information exponent**

$$\text{IE}(\ell) = \inf\{k : \|\mathbb{E}[(\mathbf{x})^{\otimes k} f^{\star}(\mathbf{x})]\|_F \neq 0\}$$

$$\text{CIE}(\ell) = \inf\{k : \|\mathbb{E}[(h_{\ell}(\mathbf{x}))^{\otimes k} f^{\star}(\mathbf{x})]\|_F = \Theta(1)\}$$

**Requires non-trivial low-degree correlations  
between intermediate features and labels**

# Conditions for learnability

**Information exponent**



**Compositional  
Information exponent**

$$\text{IE}(\ell) = \inf\{k : \|\mathbb{E}[(\mathbf{x})^{\otimes k} f^{\star}(\mathbf{x})]\|_F \neq 0\}$$

$$\text{CIE}(\ell) = \inf\{k : \|\mathbb{E}[(h_{\ell}(\mathbf{x}))^{\otimes k} f^{\star}(\mathbf{x})]\|_F = \Theta(1)\}$$

**Requires non-trivial low-degree correlations  
between intermediate features and labels**

**Matches behavior of  
real data**

	Cal-101 (30/class)	Cal-256 (60/class)
SVM (1)	44.8 ± 0.7	24.6 ± 0.4
SVM (2)	66.2 ± 0.5	39.6 ± 0.3
SVM (3)	72.3 ± 0.4	46.0 ± 0.3
SVM (4)	76.6 ± 0.4	51.3 ± 0.1
SVM (5)	<b>86.2 ± 0.8</b>	65.6 ± 0.3
SVM (7)	<b>85.5 ± 0.4</b>	<b>71.7 ± 0.2</b>
Softmax (5)	82.9 ± 0.4	65.7 ± 0.5
Softmax (7)	<b>85.4 ± 0.4</b>	<b>72.6 ± 0.1</b>

Table 7. Analysis of the discriminative information contained in each layer of feature maps within our ImageNet-pretrained convnet. We train either a linear SVM or softmax on features from different layers (as indicated in brackets) from the convnet. Higher layers generally produce more discriminative features.



# Full set of conditions

**Essential assumptions**

**Technical assumptions**

# Full set of conditions

## Essential assumptions

- $g^\star$ : information exponent 1.
- $P_k$ : information exponent/leap  
 $\leq 2$

$$\text{CIE}(1) = IE(g^\star) \times IE(P_k)$$

- expressive  $\sigma$  (non-zero Hermite) and regular.

## Technical assumptions

# Full set of conditions

## Essential assumptions

- $g^\star$ : information exponent 1.
- $P_k$ : information exponent/leap  $\leq 2$

$$\text{CIE}(1) = IE(g^\star) \times IE(P_k)$$

- expressive  $\sigma$  (non-zero Hermites) and regular.

## Technical assumptions

- $a_i^\star = 1 \ \forall i$  (symmetric targets)
- Correlation loss
- Re-initialization of layers
- expressive (non-zero Hermites) and regular.
- $P_k$ : information exponent  $\neq 1$  (to avoid spikes)
- $\mathbb{E}[\sigma(\sigma(z))\text{He}_2(z)]\mathbb{E}[P_k(z)\text{He}_2(z)] > 0$ ,  
 $\mathbb{E}[\sigma(\sigma(z))z] = 0$ ,
- $\mathbb{E}_{z \sim \mathcal{N}(0,1)}[g^\star(z)\text{He}_j(z)] = 0, \ 1 < j \leq k$

# Main theorem

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$



# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

$$W_1 = Z(W^\star) + o_d(1), z_i \sim U(S_1)$$

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

$$W_1 = Z(W^\star) + o_d(1), z_i \sim U(S_1)$$

- Second layer: upon reinitializing  $W_2 = \mathbf{0}_{d \times d}$ , one step of pre-conditioned step with batch-size  $\mathcal{O}(d^{k\epsilon+\delta})$ ,  $p_1 = \mathcal{O}(d^{k\epsilon+\delta})$ :

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

$$W_1 = Z(W^\star) + o_d(1), z_i \sim U(S_1)$$

- Second layer: upon reinitializing  $W_2 = \mathbf{0}_{d \times d}$ , one step of pre-conditioned step with batch-size  $\mathcal{O}(d^{k\epsilon+\delta})$ ,  $p_1 = \mathcal{O}(d^{k\epsilon+\delta})$ :

$$\mathbf{h}_i(\mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \mathbf{c} \mathbf{w}_3^i \mathbf{h}^\star(\mathbf{x}) + \mathbf{o}_d(1)$$



# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

$$W_1 = Z(W^\star) + o_d(1), z_i \sim U(S_1)$$

- Second layer: upon reinitializing  $W_2 = \mathbf{0}_{d \times d}$ , one step of pre-conditioned step with batch-size  $\mathcal{O}(d^{k\epsilon+\delta})$ ,  $p_1 = \mathcal{O}(d^{k\epsilon+\delta})$ :

$$\mathbf{h}_i(\mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \mathbf{c} \mathbf{w}_3^i \mathbf{h}^\star(\mathbf{x}) + \mathbf{o}_d(1)$$

- Third layer: Ridge regression on  $\mathbf{w}_3$  with samples  $\mathcal{O}(d^\delta)$ ,  $p_2 = \mathcal{O}(d^\delta)$ :

# Main theorem

$$\hat{f}(\mathbf{x}) = \mathbf{w}_3^\top \sigma \left( W_2 \sigma(W_1 \mathbf{x} + b_1) + b_2 \right)$$

$$W_1 \in \mathbb{R}^{p_1 \times d}, W_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{w}_3 \in \mathbb{R}^{p_3}$$

Let  $\delta > 0$  be arbitrarily small

- First layer:  $\mathcal{O}(\log d)$  steps of spherical SGD on  $W_1$  with correlation loss, batch-size  $\mathcal{O}(d^{1+\epsilon+\delta})$ , vanishing step-size:

$$W_1 = Z(W^\star) + o_d(1), z_i \sim U(S_1)$$

- Second layer: upon reinitializing  $W_2 = \mathbf{0}_{d \times d}$ , one step of pre-conditioned step with batch-size  $\mathcal{O}(d^{k\epsilon+\delta})$ ,  $p_1 = \mathcal{O}(d^{k\epsilon+\delta})$ :

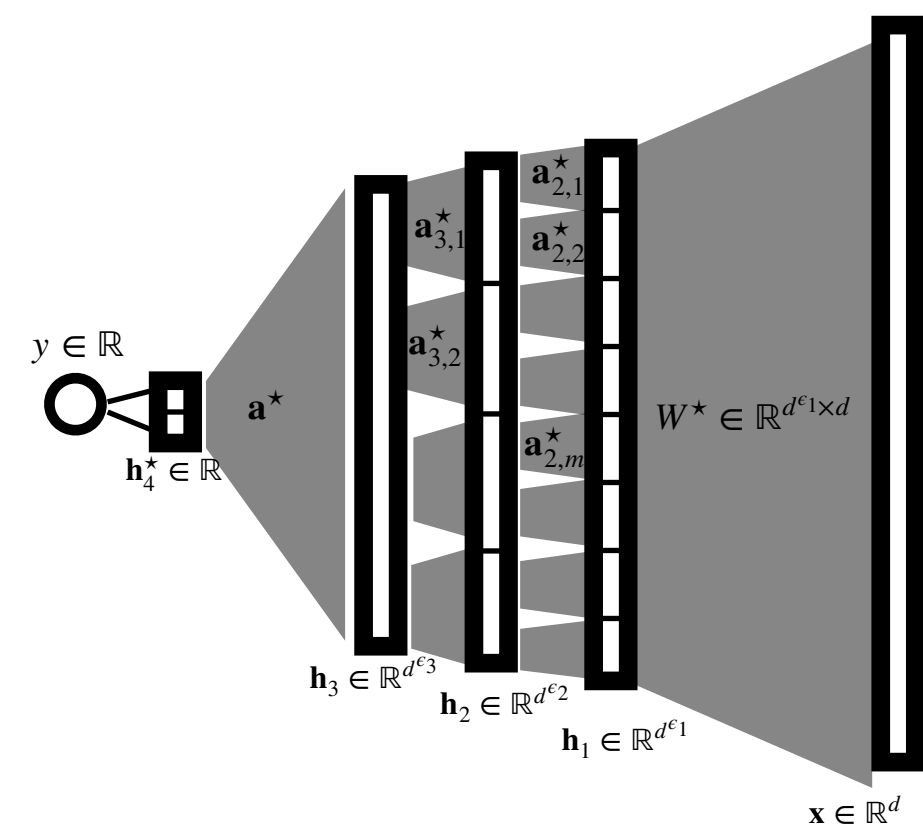
$$\mathbf{h}_i(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) = \mathbf{c} \mathbf{w}_3^i \mathbf{h}^\star(\mathbf{x}) + \mathbf{o}_d(1)$$

- Third layer: Ridge regression on  $\mathbf{w}_3$  with samples  $\mathcal{O}(d^\delta)$ ,  $p_2 = \mathcal{O}(d^\delta)$ :

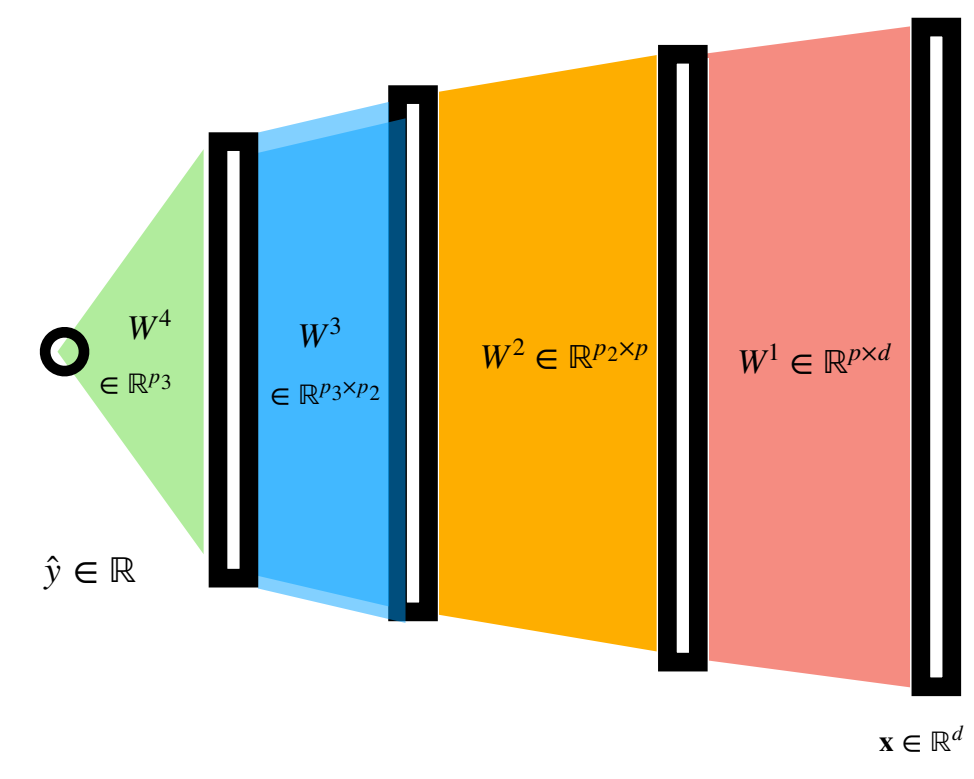
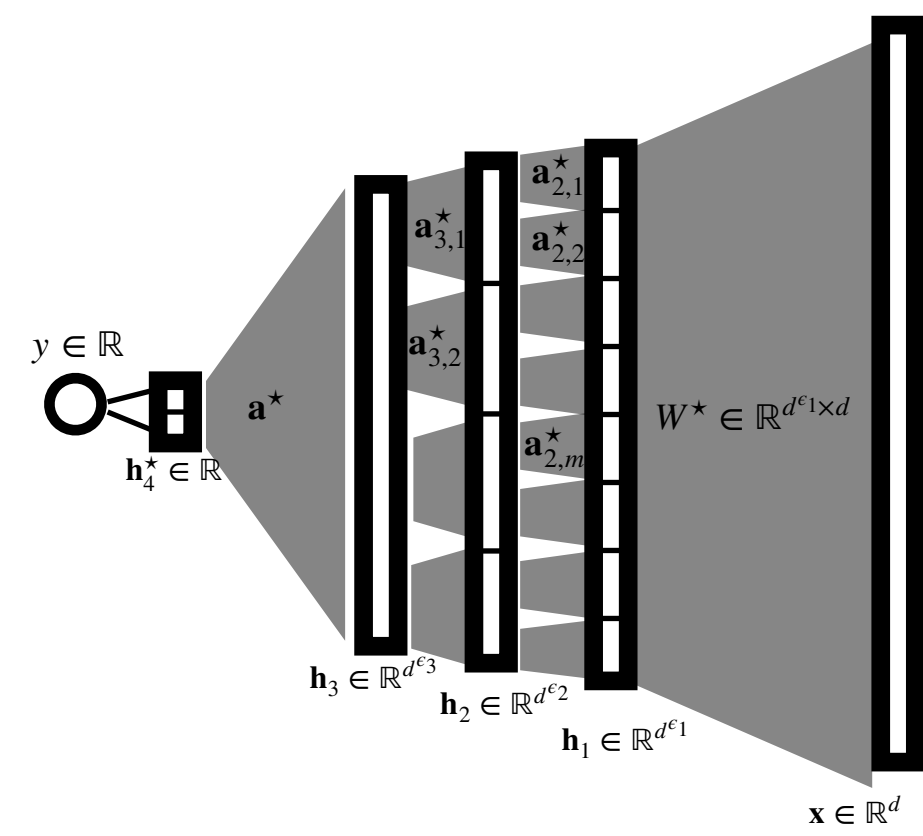
$$\hat{f}(\mathbf{x}) = \mathbf{f}^\star(\mathbf{x}) + \mathbf{o}_d(1)$$

# General Depth

# General Depth

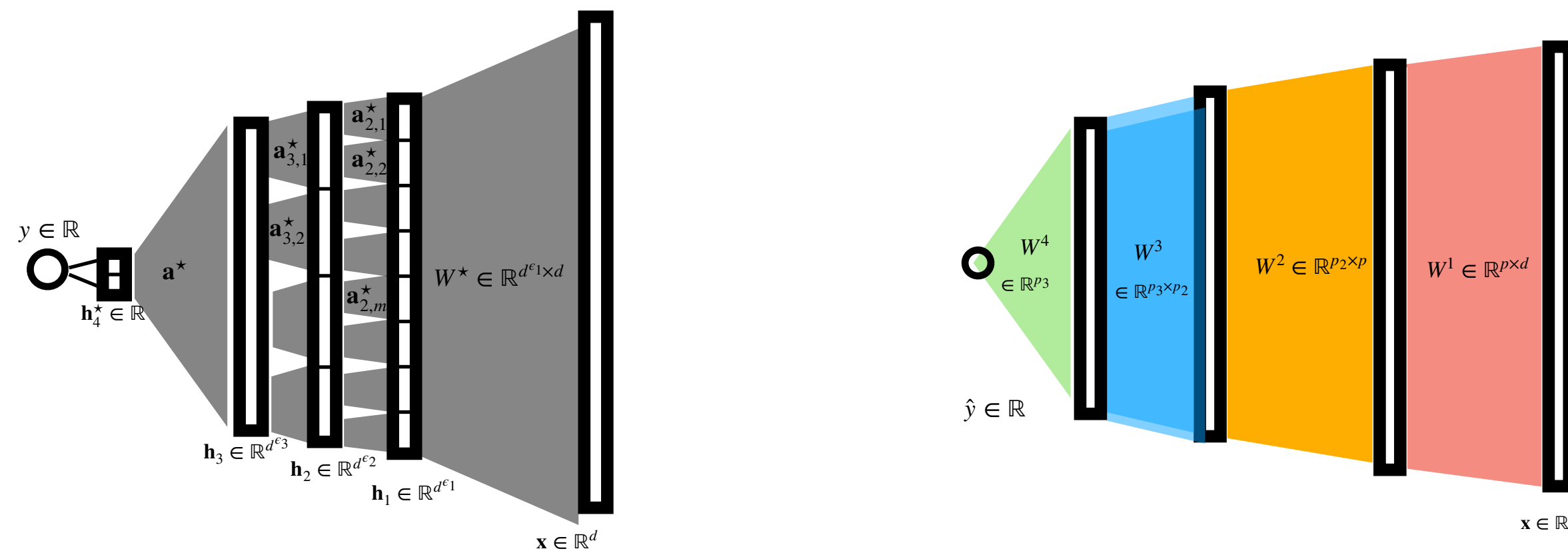


# General Depth



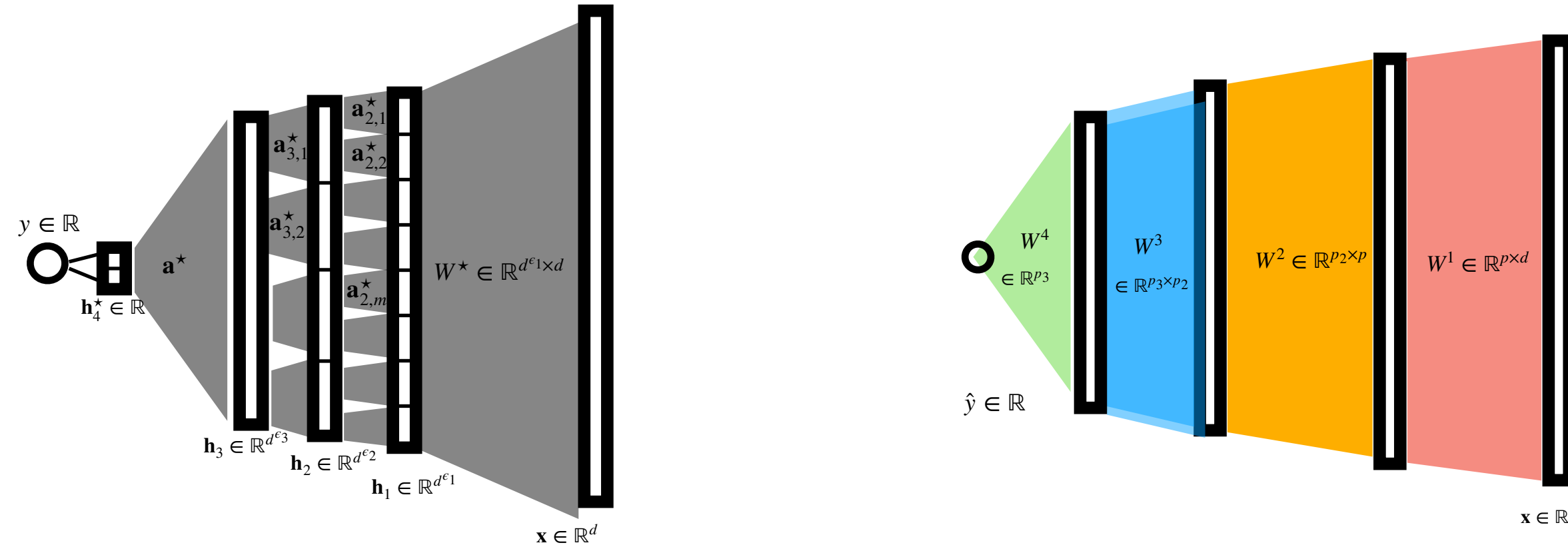


# General Depth



- Conditional on perfect spherical recovery for  $W_{L-2}$ , the same picture holds for the last two layers.
- Key idea: Features are independent by tree-structure, asymptotically Gaussian, and maintain nice tails (hypercontractivity).

# General Depth



- Conditional on perfect spherical recovery for  $W_{L-2}$ , the same picture holds for the last two layers.
- Key idea: Features are independent by tree-structure, asymptotically Gaussian, and maintain nice tails (hypercontractivity).

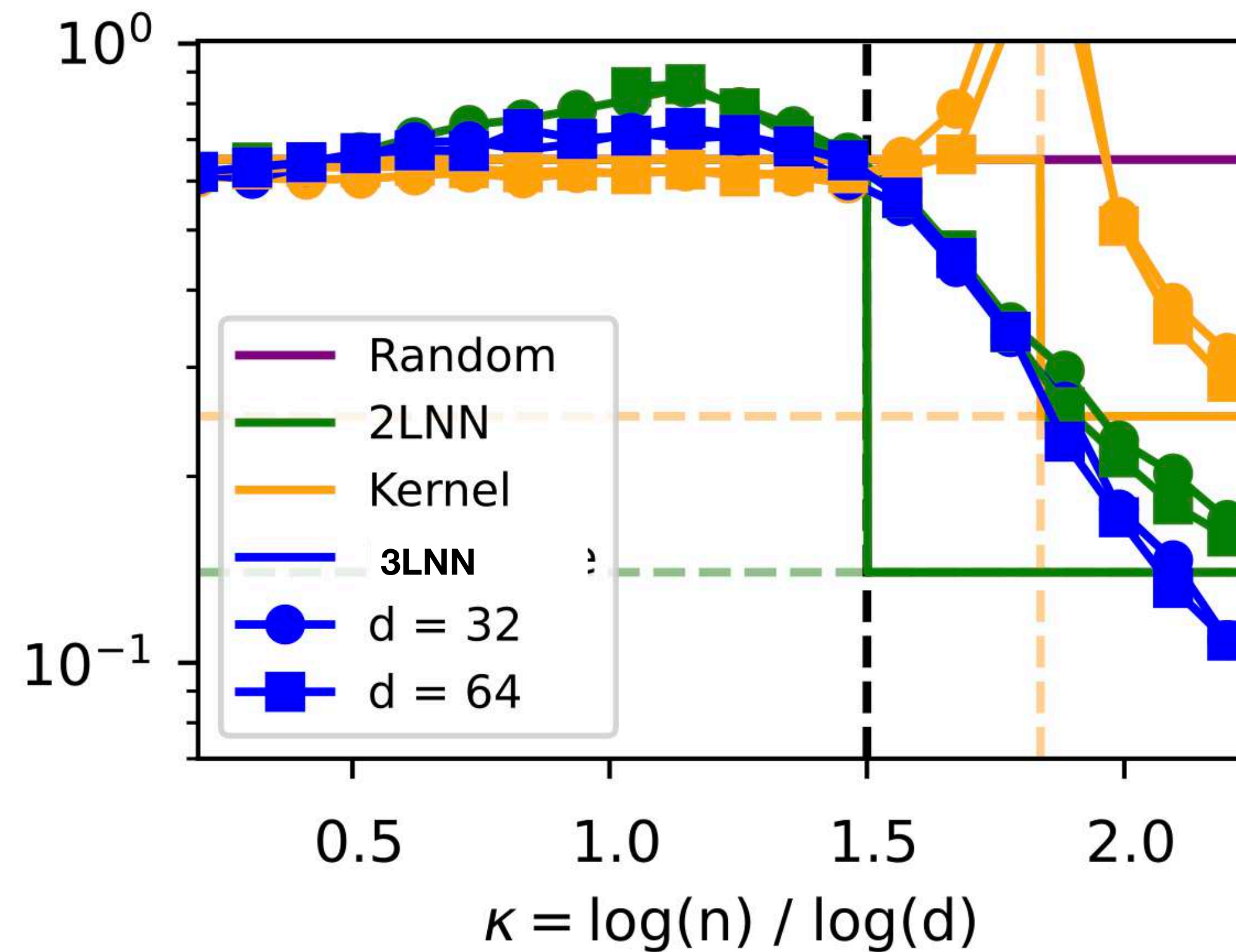
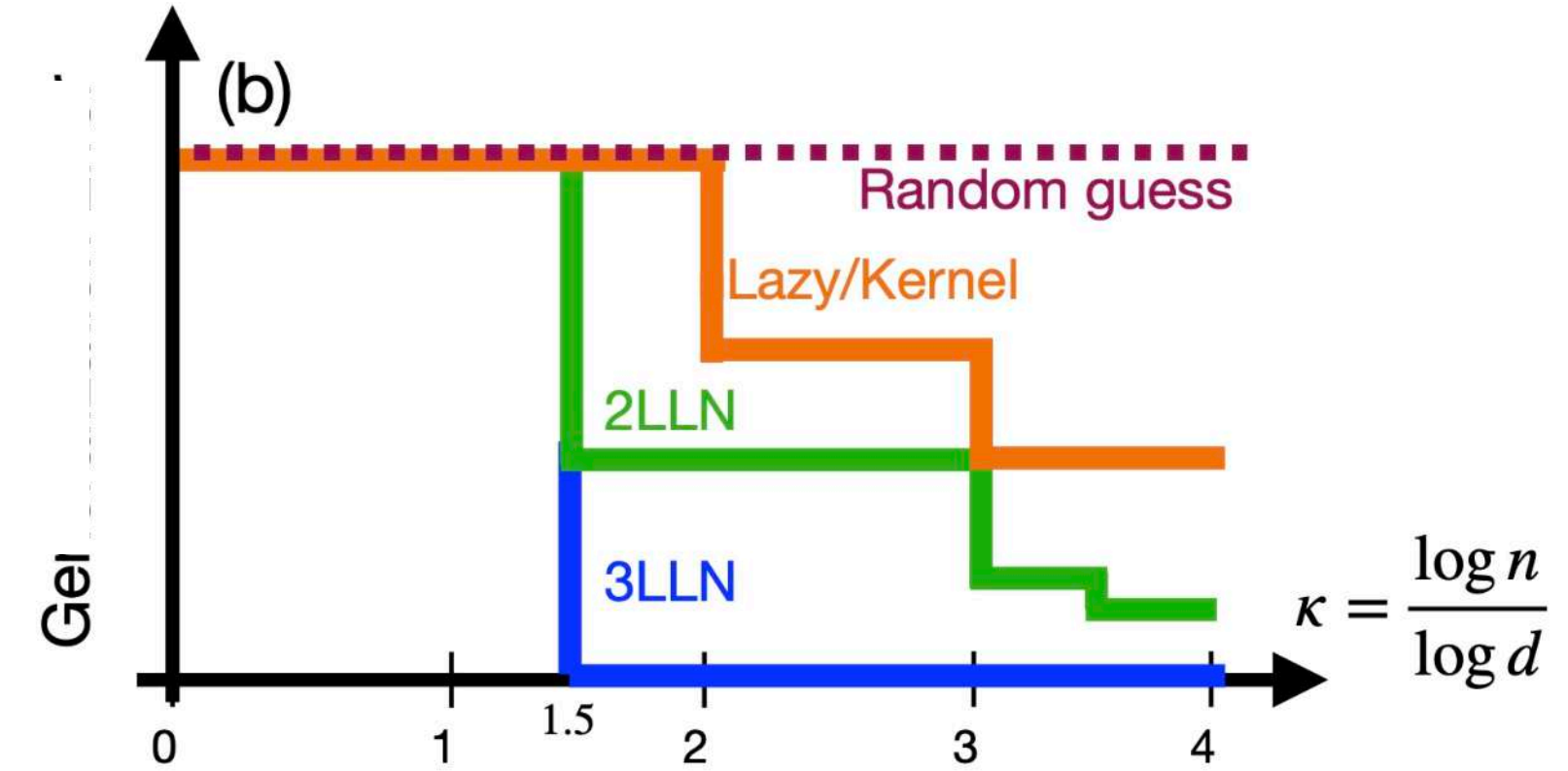
**Theorem 2.** For  $L \in \mathbb{N}$ , let  $f^*(\mathbf{x})$  denote a target as in Eq. (6) with  $r = 1$ , and let  $\delta', \delta$  be arbitrary reals satisfying  $0 < \delta < \delta' < 1$ . Consider a model of the form  $\hat{f}_\theta(\mathbf{x}) = \mathbf{w}_L^\top \sigma(W_{L-1} \sigma(W h_{L-1}^*(\mathbf{x})))$  with  $W \in \mathbb{R}^{p_{L-2} \times d^{L-2}}$  having  $p_{L-2} = \Theta(d^{k_{\varepsilon_{L-2}} + \delta'})$  rows independently sampled as  $\mathbf{w}_i \sim U(\mathcal{S}_{d^{L-2}}(1))$ . Under Ass. 1-3, after a single step of pre-conditioned SGD on  $W_{L-1}$  with batch-size  $\Theta(d^{k_{\varepsilon_{L-2}} + \delta})$ , step-size  $\Theta(\sqrt{p_{L-1}})$ , the pre-activations  $h_{L-1}(\mathbf{x}) := W_{L-1} \sigma(W h_{L-1}^*(\mathbf{x}))$  satisfy, for a constant  $c > 0$ :

$$h_{L-1}(\mathbf{x}) = c \mathbf{w}_L h_L^*(\mathbf{x}) + o_d(1), \quad (22)$$

# Example

$$f^*(\mathbf{x}) = \tanh \left( \frac{\mathbf{a}^{*\top} P_3(W^*\mathbf{x})}{\sqrt{d^{\varepsilon_1=1/2}}} \right)$$

$$P_3(x) = \text{He}_2(x) + \text{He}_3(x)$$



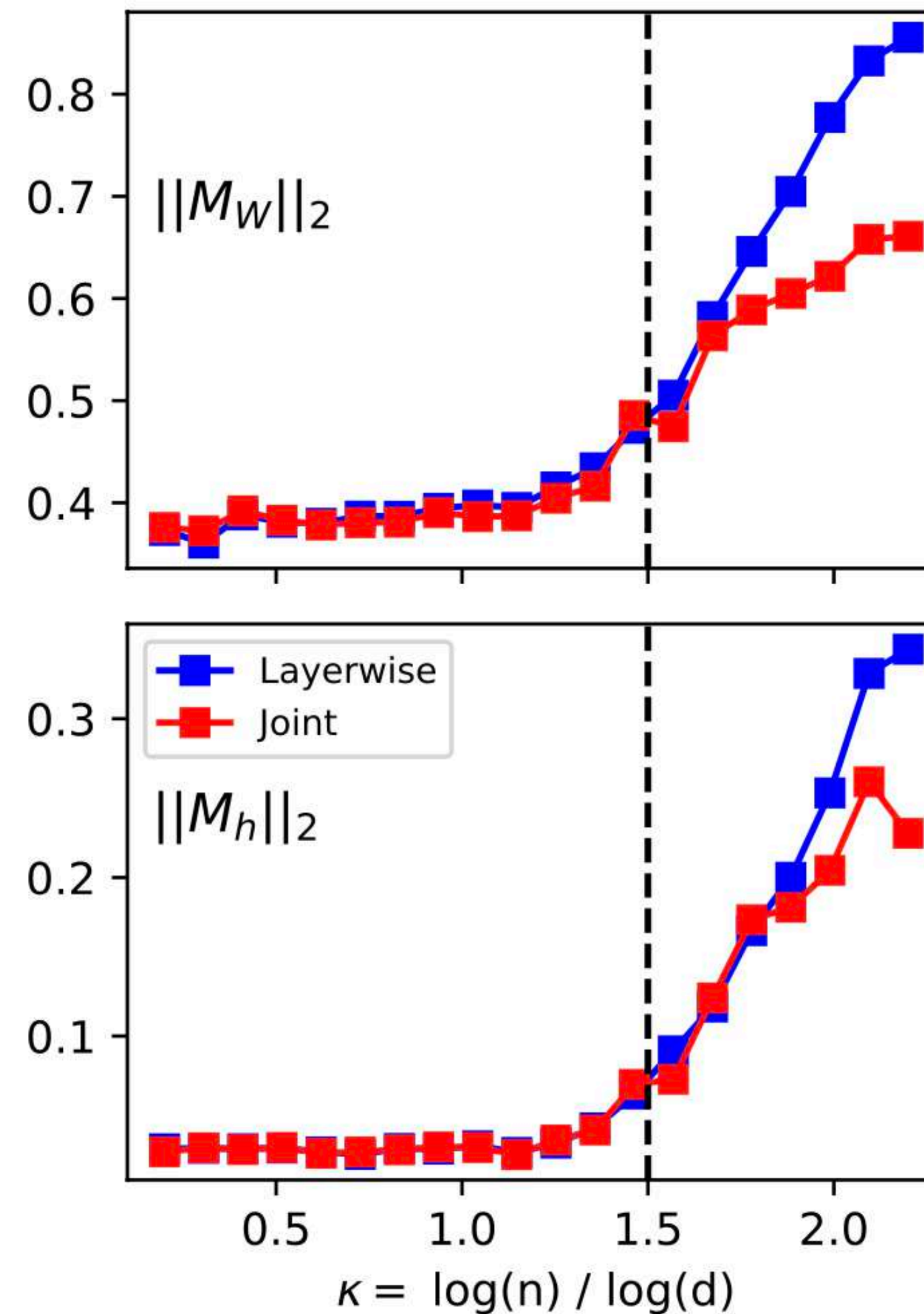


# Overlaps in parameter and function space

$$f^*(\mathbf{x}) = \tanh \left( \frac{\mathbf{a}^{*\top} P_3(W^*\mathbf{x})}{\sqrt{d^{\varepsilon_1=1/2}}} \right)$$

$$P_3(x) = \text{He}_2(x) + \text{He}_3(x)$$

$$M_W = \frac{W_1 W^*}{\|W_1\|_2},$$
$$M_h = \frac{\mathbb{E}[\mathbf{h}(\mathbf{z}) \mathbf{h}^*(\mathbf{z})]}{\sqrt{\mathbb{E}[\mathbf{h}(\mathbf{z})^2]}}.$$









**Key ideas and proof sketches**

# Recursive Expansion of the Target

# Recursive Expansion of the Target

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}^{\star^\top} P_k (W^\star \mathbf{x})}{\sqrt{d^\varepsilon}} \right), \mathbf{x} \in \mathbb{R}^d$$

# Recursive Expansion of the Target

$$f^{\star}(\mathbf{x}) = g^{\star} \left( \frac{\mathbf{a}^{\star \top} P_k (W^{\star} \mathbf{x})}{\sqrt{d^{\varepsilon}}} \right), \mathbf{x} \in \mathbb{R}^d \quad h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k (\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

# Recursive Expansion of the Target

$$f^\star(\mathbf{x}) = g^\star\left(\frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\varepsilon}}\right), \mathbf{x} \in \mathbb{R}^d \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\varepsilon}} \in \mathbb{R}$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2 \mathbf{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$



# Recursive Expansion of the Target

$$f^\star(\mathbf{x}) = g^\star\left(\frac{\mathbf{a}^{\star\top} P_k(W^\star \mathbf{x})}{\sqrt{d^\varepsilon}}\right), \mathbf{x} \in \mathbb{R}^d \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\varepsilon}} \in \mathbb{R}$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2 \mathbf{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$

**!!!  $He_k(h^\star(\mathbf{x}))$  can contribute low-degree terms**

# Composition of Hermites

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$

# Composition of Hermities

$$h(\mathbf{x}) = \frac{1}{\sqrt{\mathbf{d}^\epsilon}} \sum_{i=1}^{\mathbf{d}^\epsilon} \text{He}_{\mathbf{k}}(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$
$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{\mathbf{m} \mathbf{d}^{\mathbf{m} \epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_{\mathbf{k}}(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$

$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$



# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$

$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$

High Hermite-degree

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$

$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$

$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$

- The dominant term in the expansion is along  $h^\star(\mathbf{x})$ .

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$
$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$

- The dominant term in the expansion is along  $h^\star(\mathbf{x})$ .
- Generalization of the “approximate Stein’s Lemma” in Nichani et al. 2023.

# Composition of Hermites

$$h(\mathbf{x}) = \frac{1}{\sqrt{d^\epsilon}} \sum_{i=1}^{d^\epsilon} \text{He}_k(\langle \mathbf{w}_i^\star, \mathbf{x} \rangle)$$
$$\text{He}_m(h(\mathbf{x})) \approx \frac{1}{\sqrt{m d^{m\epsilon_1}}} \sum_{\text{distinct subsets}} \prod_{s_i} \text{He}_k(\langle \mathbf{w}_{s_i}^\star, \mathbf{x} \rangle)$$

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \underbrace{\mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots}_{\text{High Hermite-degree}}$$

- The dominant term in the expansion is along  $h^\star(\mathbf{x})$ .
- Generalization of the “approximate Stein’s Lemma” in Nichani et al. 2023.



# Recovery by the first layer

# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \mathbf{H} \mathbf{e}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$

# Recovery by the first layer

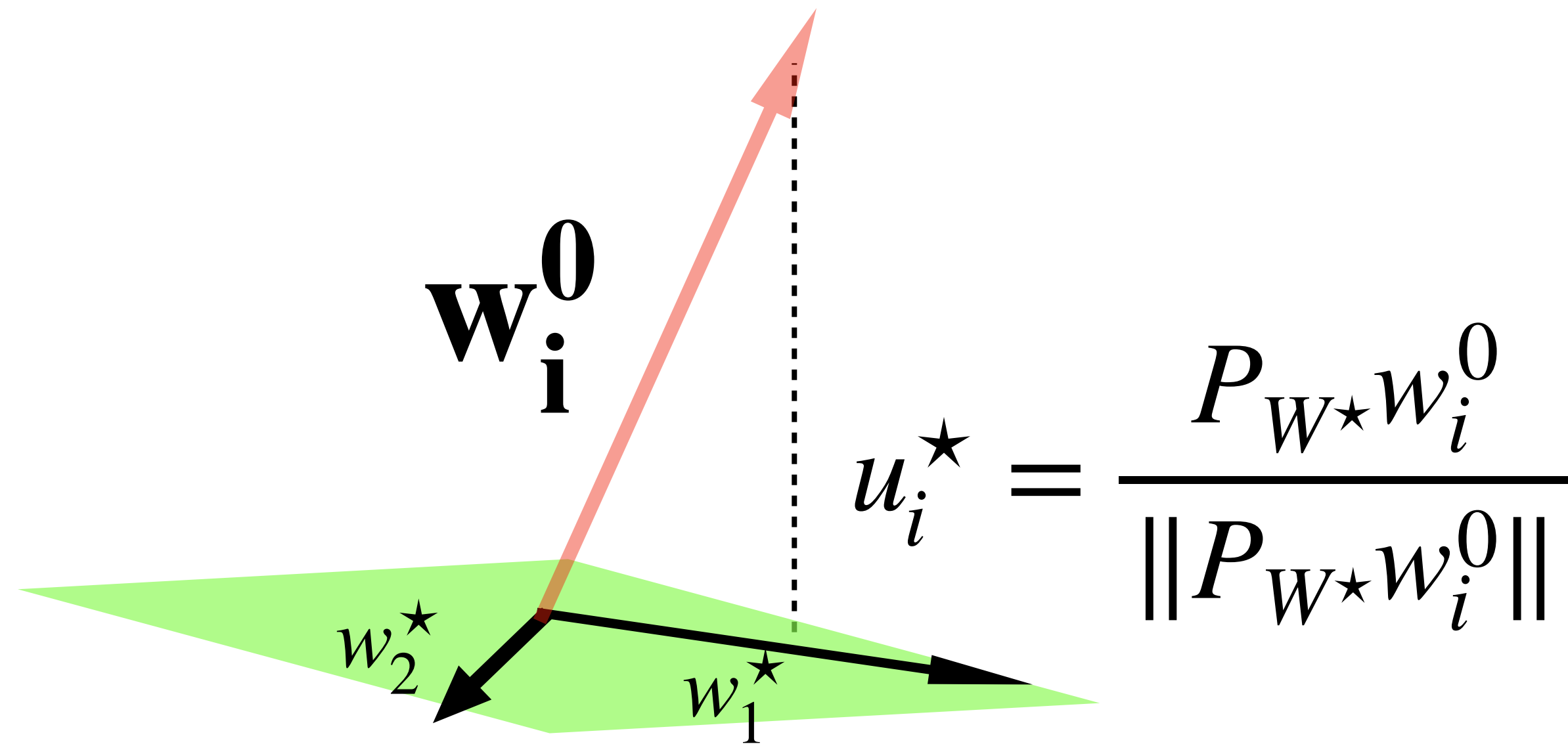
$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2 \mathbf{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$

**Updates to  $W_1 \approx \text{SGD on}$**

$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\varepsilon}} \in \mathbb{R}$$

# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \mathbf{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$

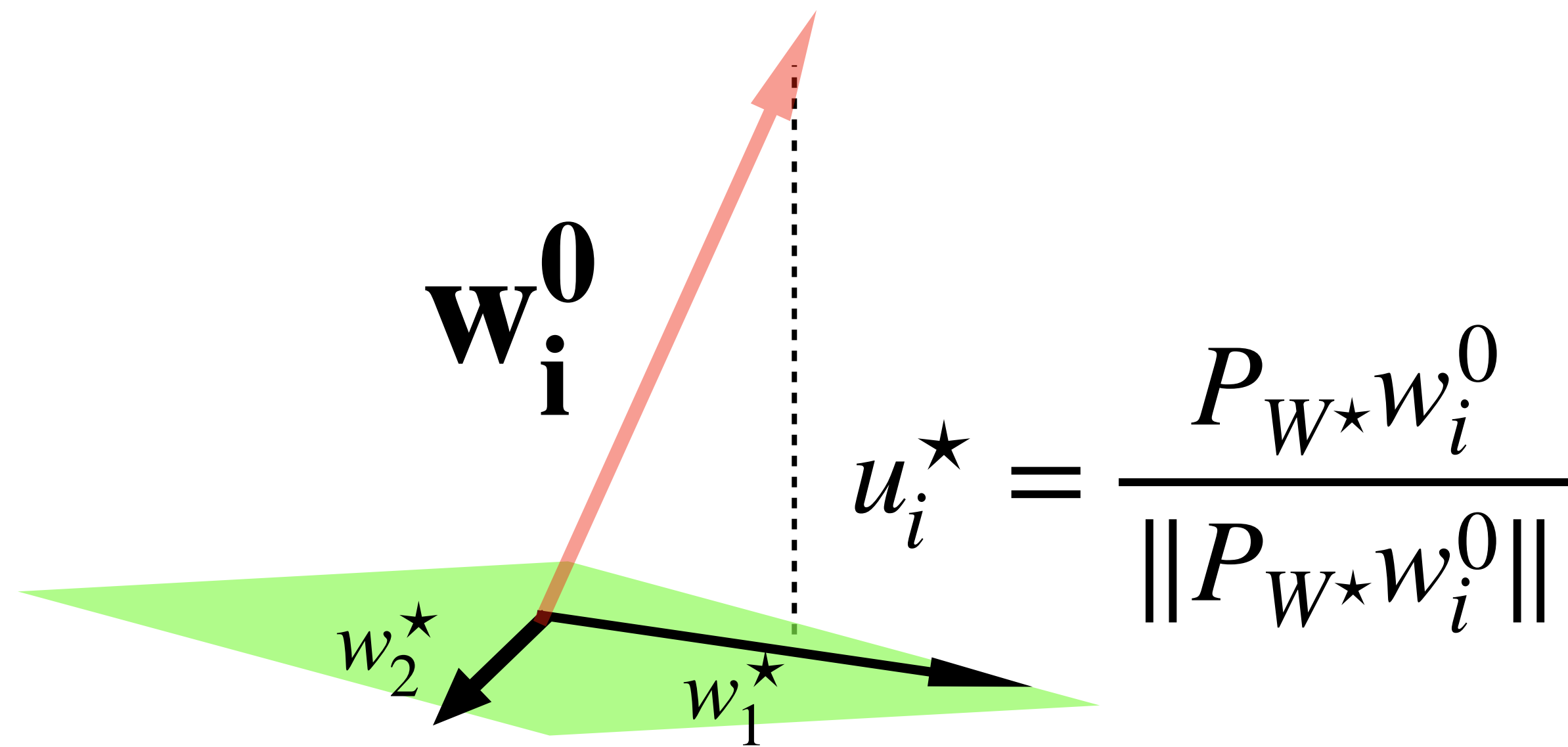


**Updates to  $W_1 \approx \text{SGD on}$**

$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$



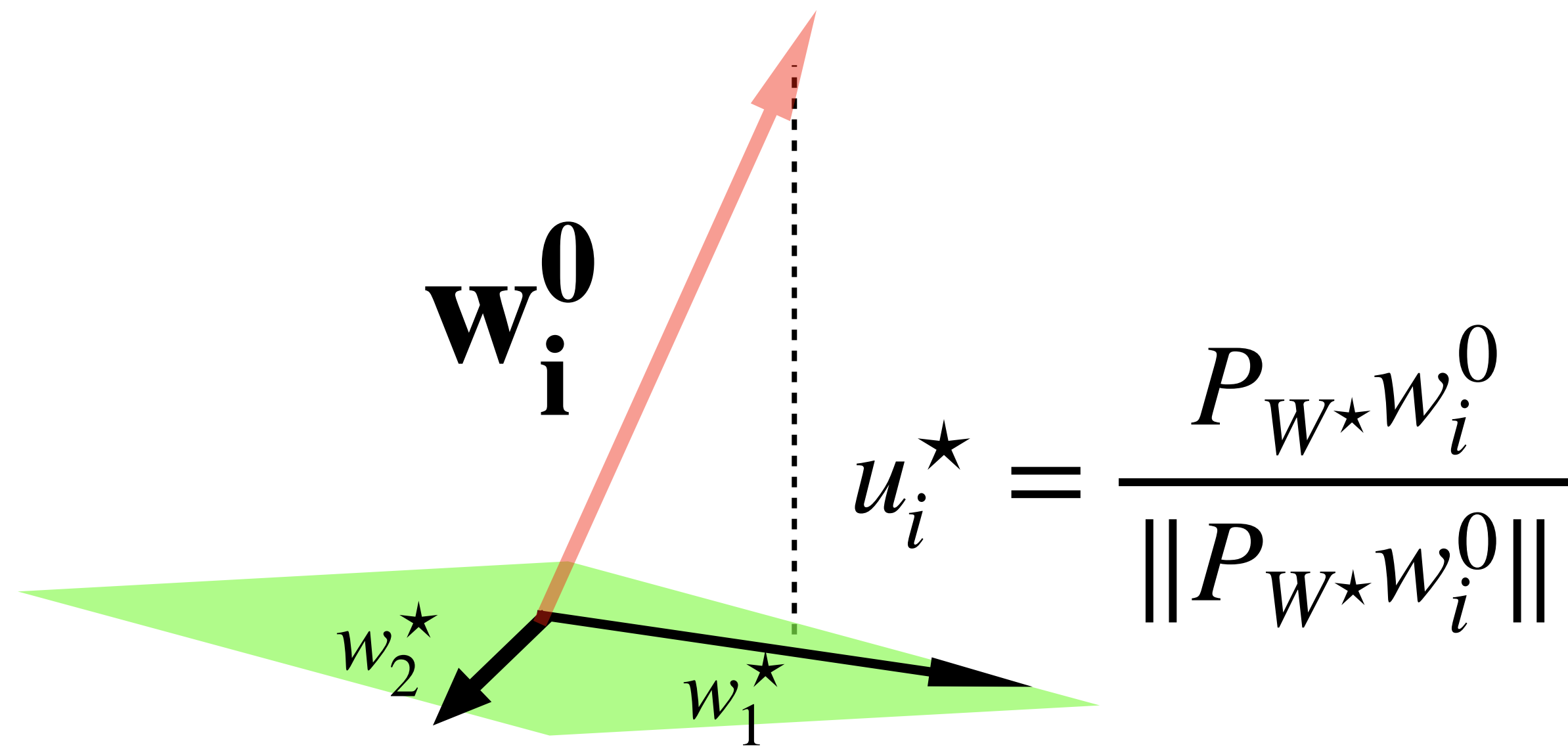
**Gradient dominated by initial direction**

**Updates to  $W_1 \approx \text{SGD}$  on**

$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$



**Gradient dominated by initial direction**

**Updates to  $W_1 \approx \text{SGD}$  on**

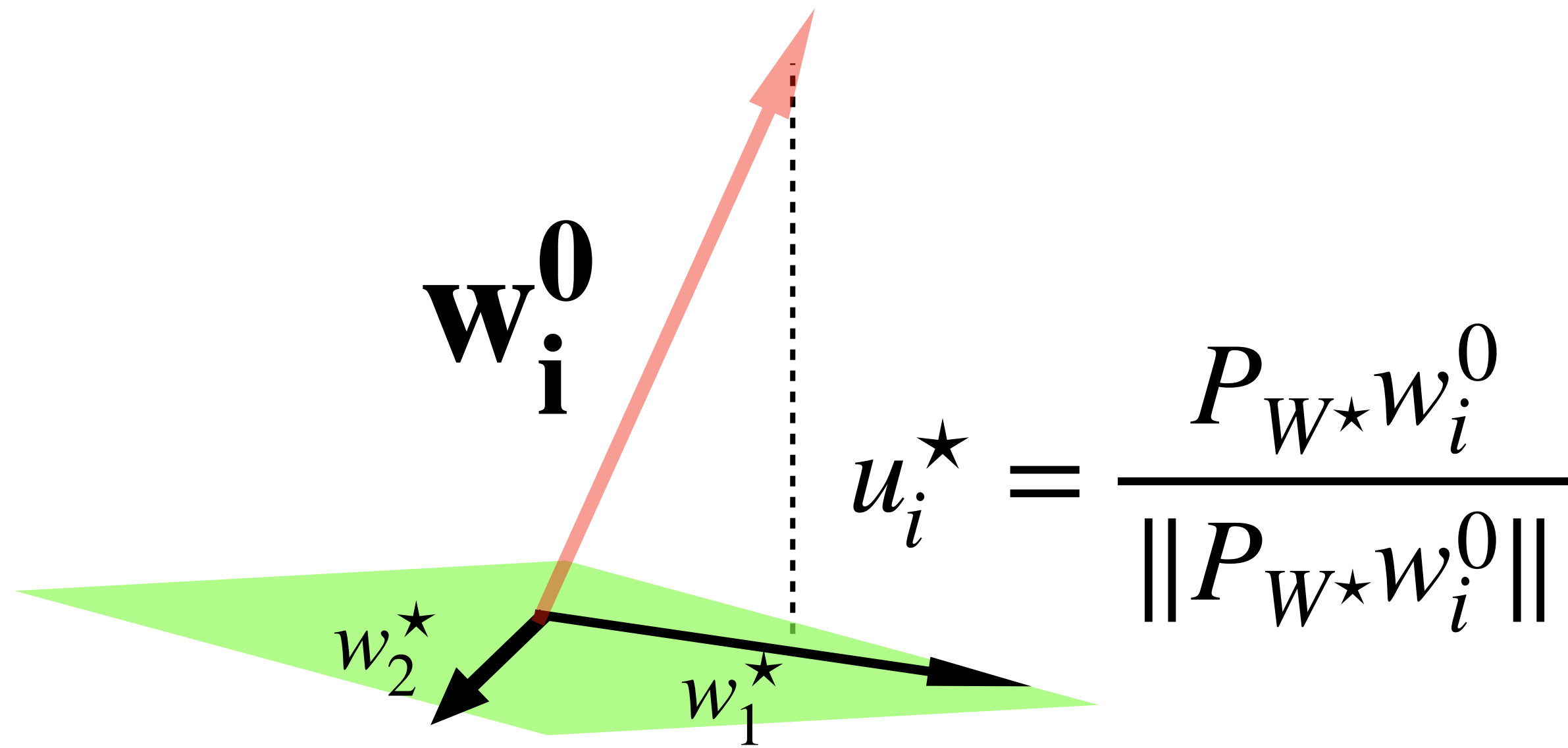
$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$\frac{d\mathbf{w}_i}{dt} \propto u_i^\star + \text{noise}$$



# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$



**Gradient dominated by initial direction**

**Updates to  $W_1 \approx \text{SGD}$  on**

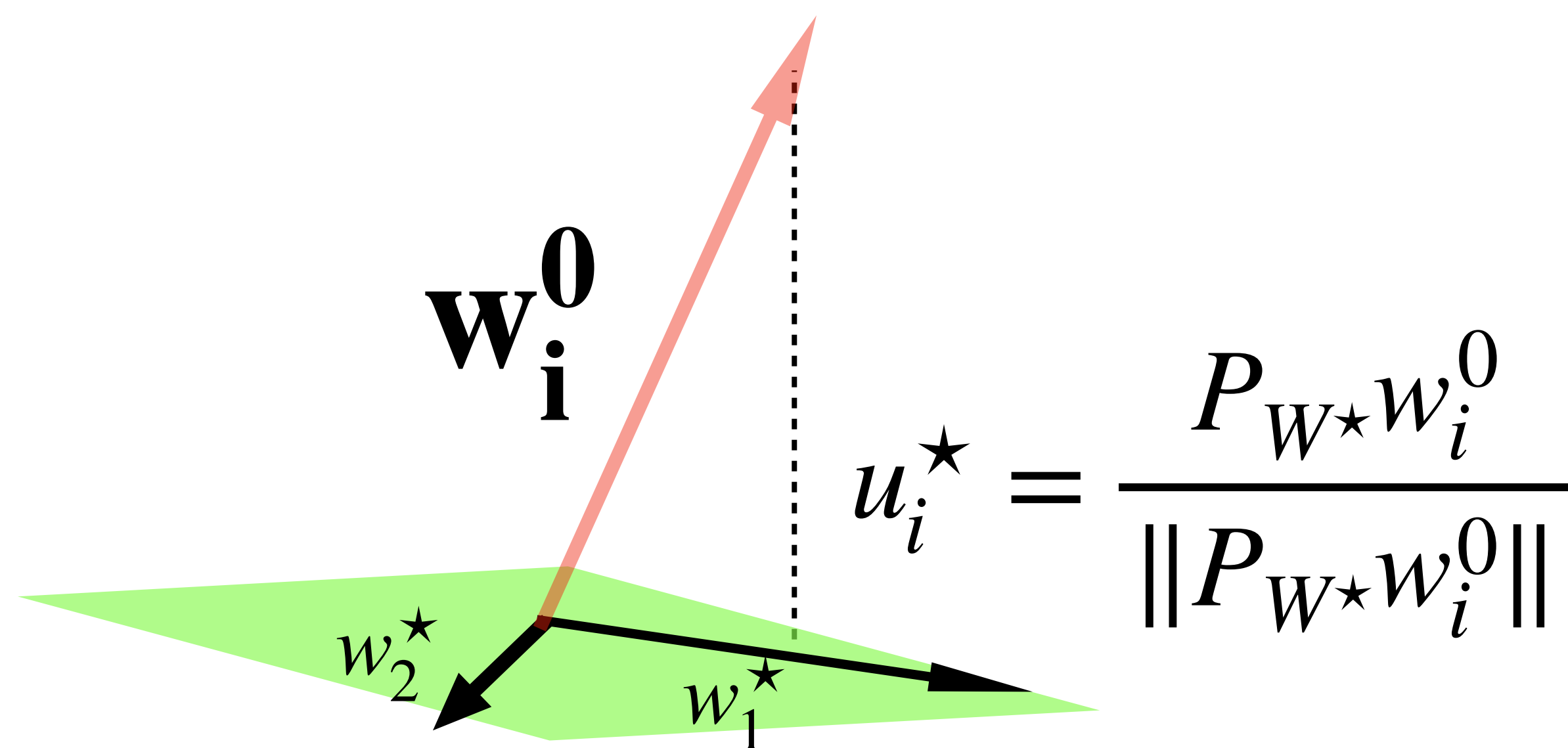
$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$\frac{d\mathbf{w}_i}{dt} \propto u_i^\star + \text{noise}$$

**drift+martingale**

# Recovery by the first layer

$$g^\star(h_\star(\mathbf{x})) \approx \mu_1 \mathbf{h}^\star(\mathbf{x}) + \mu_2^\star \text{He}_2(\mathbf{h}^\star(\mathbf{x})) + \dots$$



**Updates to  $W_1 \approx \text{SGD}$  on**

$$h^\star = \mathbf{a}^\star \cdot \frac{P_k(W^\star \mathbf{x})}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$\frac{d\mathbf{w}_i}{dt} \propto u_i^\star + \text{noise}$$

**drift+martingale**

**Gradient dominated by initial direction**

**higher-order terms suppressed by vanishing specialization along  $w_i^\star$  + vanishing step size.**

# Recovery by the second layer

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

**Gradient updates**

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

**Gradient updates**

$\approx$



# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

**Gradient updates**

$\approx$

**Projections on the conjugate  
Kernel\***

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

Gradient updates

$\approx$

Projections on the conjugate  
Kernel\*

$$\Delta W_2^i \propto w_i^3 \sigma(W_1 X^\top) f^\star(X) \odot \sigma'(\mathbf{h}_i^t(\mathbf{X}))$$

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

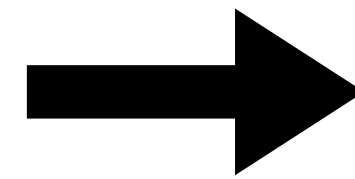
Gradient updates

$\approx$

Projections on the conjugate  
Kernel\*

$$\Delta h_{i,2}^t(\mathbf{x}) \approx \langle \sigma(W_1 \mathbf{x}), \Delta W_i \rangle$$

$$\Delta W_2^i \propto w_i^3 \sigma(W_1 X^\top) f^\star(X) \odot \sigma'(\mathbf{h}_i^t(\mathbf{X}))$$



# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

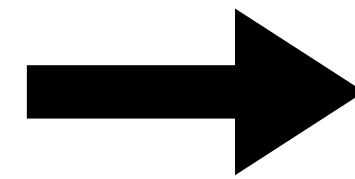
Gradient updates

$\approx$

Projections on the conjugate  
Kernel\*

$$\Delta h_{i,2}^t(\mathbf{x}) \approx \langle \sigma(W_1 \mathbf{x}), \Delta W_i \rangle$$

$$\Delta W_2^i \propto w_i^3 \sigma(W_1 X^\top) f^\star(X) \odot \sigma'(\mathbf{h}_i^t(X))$$



$$\Delta \mathbf{h}_i(\mathbf{x}) \propto \underbrace{w_i^3 \sigma(W_1 \mathbf{x})^\top \sigma(W_1 X^\top)}_{\text{Conjugate Kernel}} \underbrace{f^\star(X) \odot \sigma'(\mathbf{h}_i^t(X))}_{\text{Perturbed target}}$$

Conjugate Kernel

Perturbed target

# Recovery by the second layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \in \mathbb{R}^{p_2}, X \in \mathbb{R}^{n \times d}$$

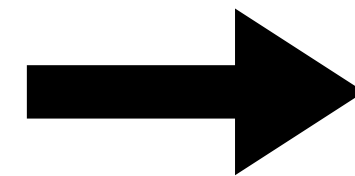
Gradient updates

$\approx$

Projections on the conjugate  
Kernel\*

$$\Delta h_{i,2}^t(\mathbf{x}) \approx \langle \sigma(W_1 \mathbf{x}), \Delta W_i \rangle$$

$$\Delta W_2^i \propto w_i^3 \sigma(W_1 X^\top) f^\star(X) \odot \sigma'(\mathbf{h}_i^t(X))$$



$$\Delta \mathbf{h}_i(\mathbf{x}) \propto \underbrace{w_i^3 \sigma(W_1 \mathbf{x})^\top \sigma(W_1 X^\top)}_{\text{Conjugate Kernel}} \underbrace{f^\star(X) \odot \sigma'(\mathbf{h}_i^t(X))}_{\text{Perturbed target}}$$

Conjugate Kernel

Perturbed target

$$\Delta \mathbf{h}_i(\mathbf{x}) \propto w_i^3 \sigma(W_1 \mathbf{x})^\top \sigma(W_1 X^\top) f^\star(X)$$

# Dimension reduction in action



# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$

# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$

# Dimension reduction in action

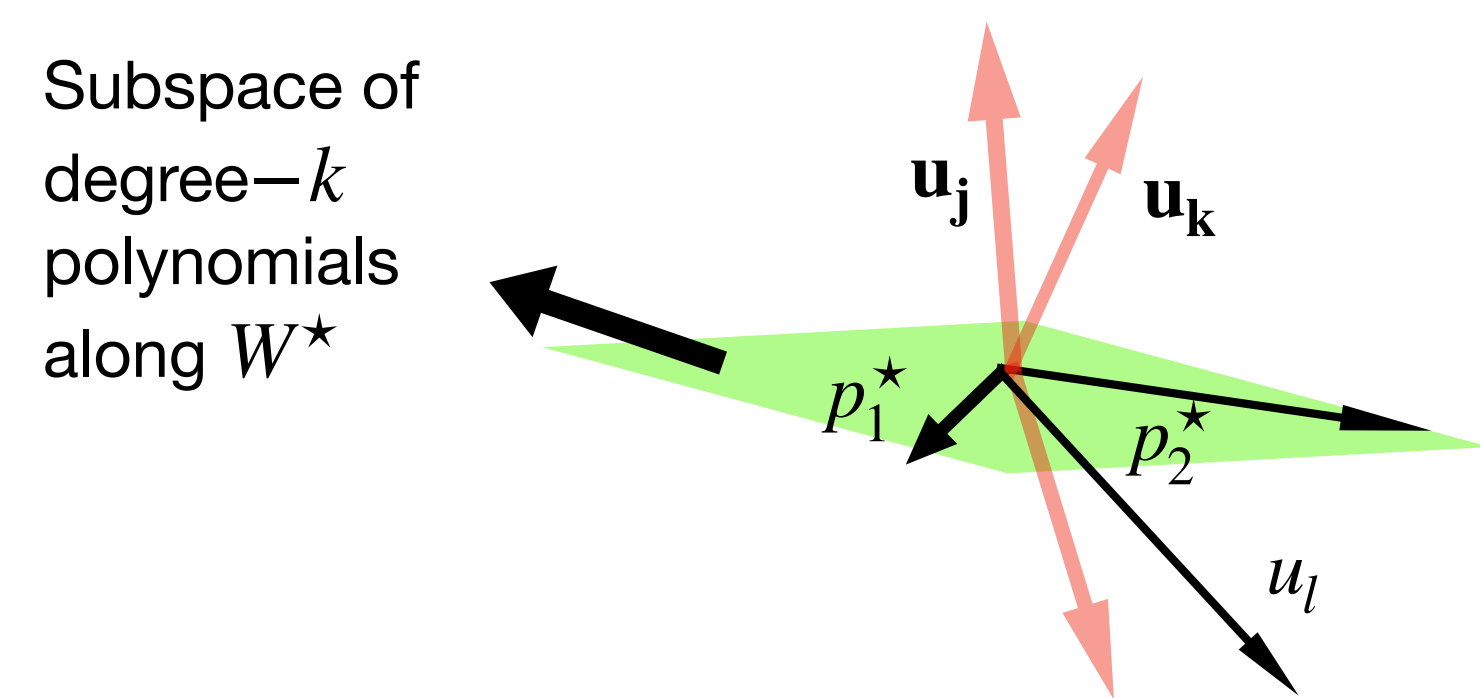
$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$

# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$

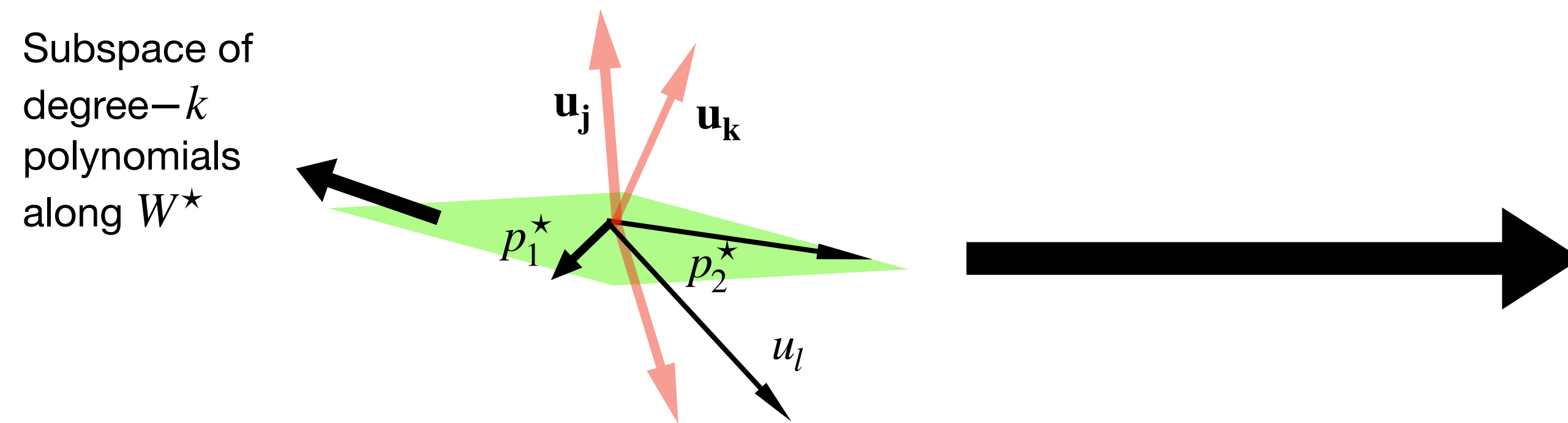


$K_{CK}(\mathbf{x}, \mathbf{x}')$  before training  $\mathbf{W}_1$

# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{W}_1 \mathbf{x})^\top \sigma(\mathbf{W}_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$



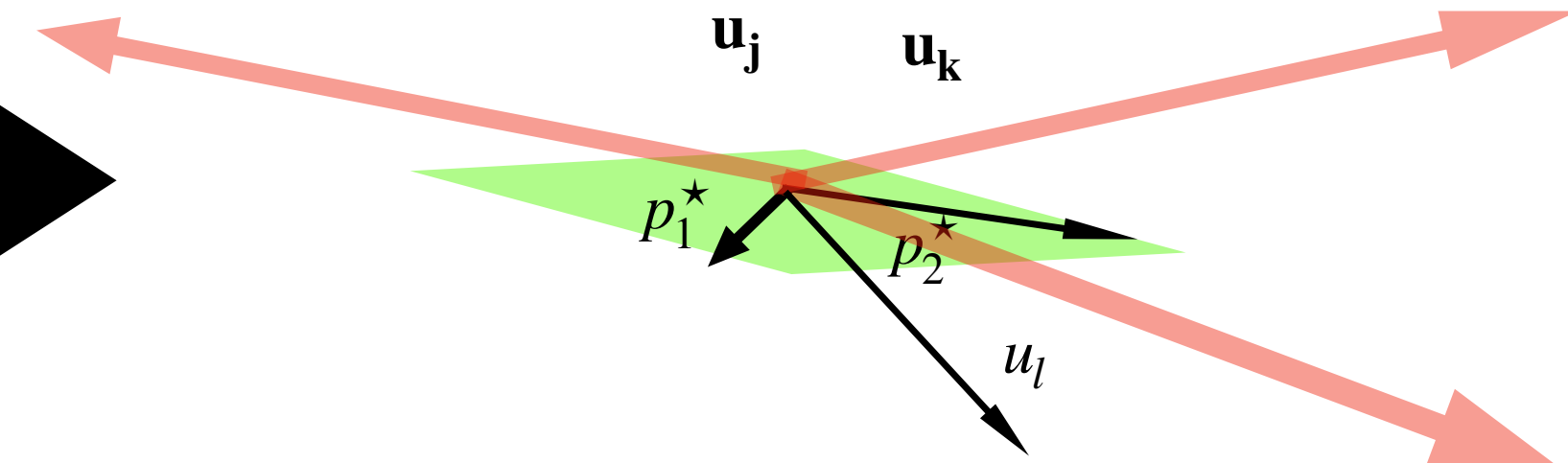
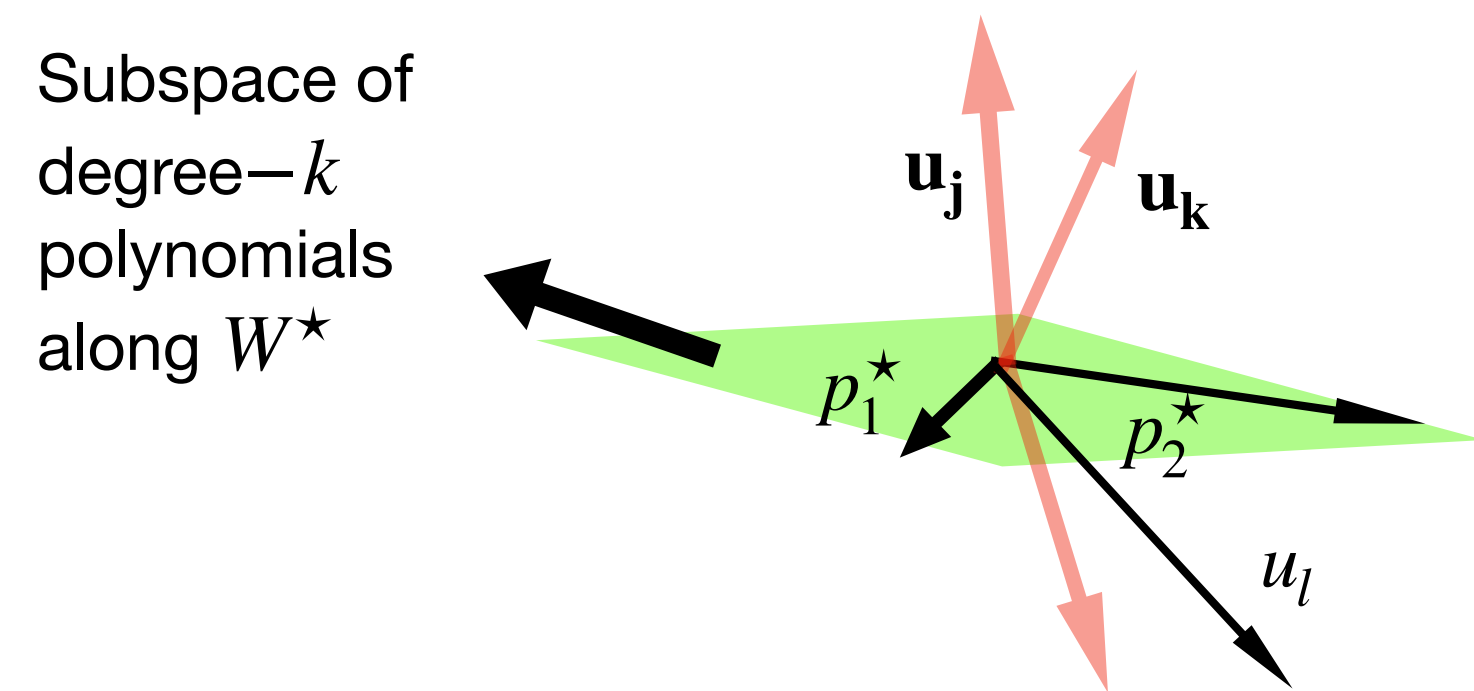
$K_{CK}(\mathbf{x}, \mathbf{x}')$  before training  $\mathbf{W}_1$



# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(W_1 \mathbf{x})^\top \sigma(W_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$



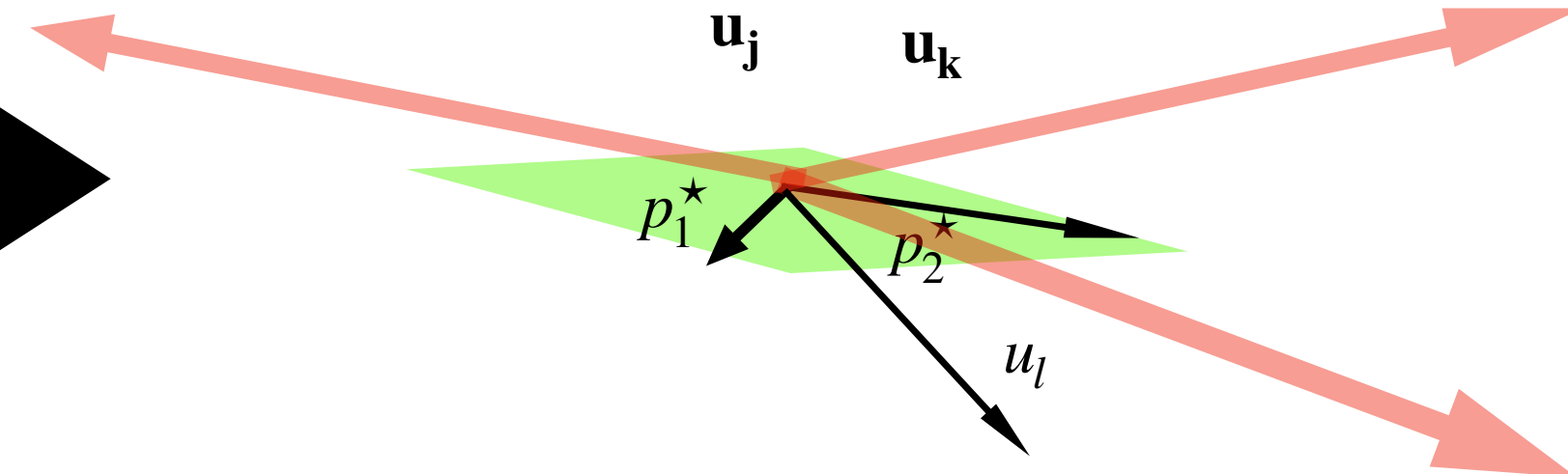
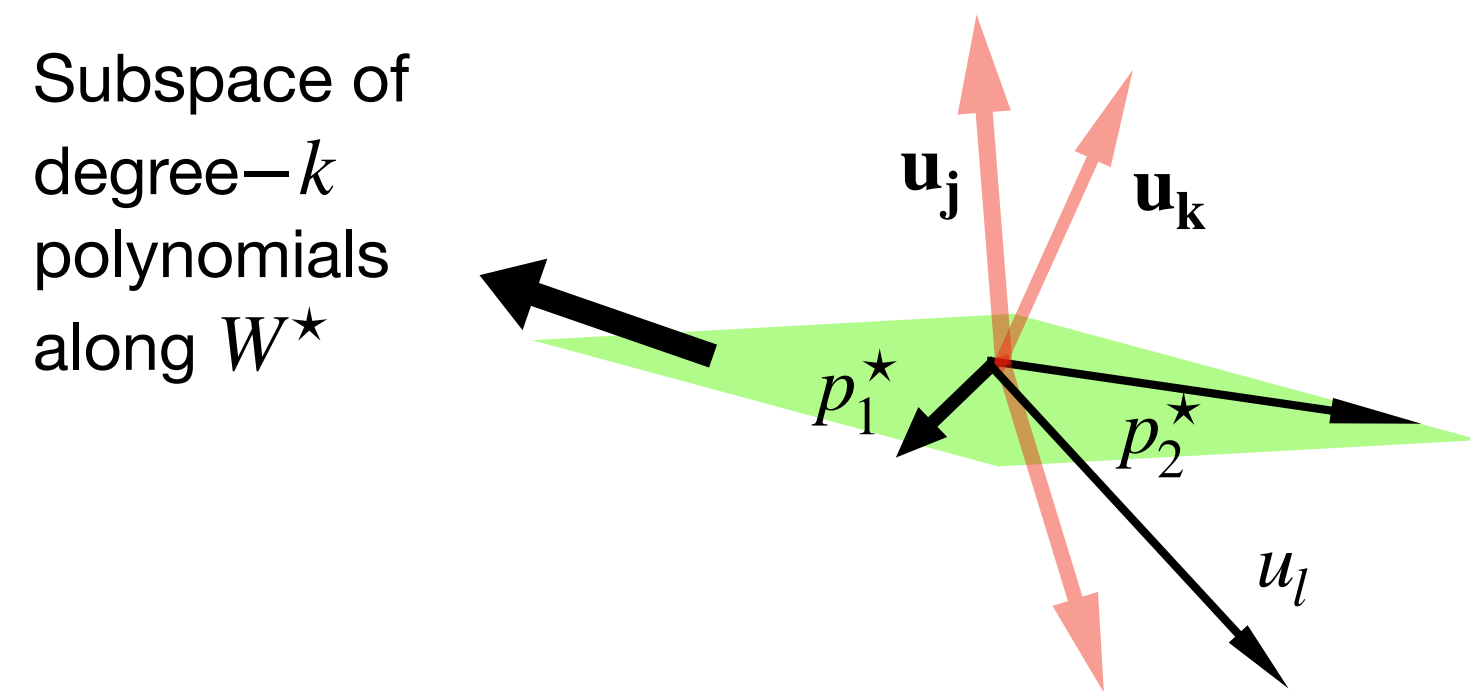
$K_{CK}(\mathbf{x}, \mathbf{x}')$  after training  $W_1$

$K_{CK}(\mathbf{x}, \mathbf{x}')$  before training  $W_1$

# Dimension reduction in action

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sigma(W_1 \mathbf{x})^\top \sigma(W_1 \mathbf{x}')$$

$$K_{CK}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}')$$



$K_{CK}(\mathbf{x}, \mathbf{x}')$  after training  $W_1$

$K_{CK}(\mathbf{x}, \mathbf{x}')$  before training  $W_1$

## Caveats:

- Conjugate Kernel ill-conditioned,  $\lambda_k = \frac{1}{d^{\epsilon k}} \implies \mathcal{O}(d^{\epsilon k})$  steps for convergence  $\implies \mathcal{O}(d^{2\epsilon k})$  sample complexity.

- Fix: Pre-conditioning:  $W_2 = W_2 - \eta \left( \frac{1}{n} \sigma(W_1 X^\top) \sigma(W_1 X)^\top \right)^{-1} \nabla_{W_2} \mathcal{L}$

# Spike+ bulk decomposition of Gram matrix

# Spike+ bulk decomposition of Gram matrix

$$\mathbf{x}_{\star} = \mathbf{W}^{\star} \mathbf{x}$$

# Spike+ bulk decomposition of Gram matrix

$$\mathbf{x}_{\star} = \mathbf{W}^{\star} \mathbf{x} \qquad h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\varepsilon}}} \in \mathbb{R}$$

# Spike+ bulk decomposition of Gram matrix

$$\mathbf{X}_\star = \mathbf{W}^\star \mathbf{X} \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\varepsilon}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_\star X_\star^T}{d} + \kappa_2^2 \frac{H_2(X_\star) H_2(X_\star)^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_\star) H_k(X_\star)^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_\star) H_{k+1}(X_\star)^T}{d^{k+1}} + \dots$$



# Spike+ bulk decomposition of Gram matrix

$$\mathbf{X}_{\star} = \mathbf{W}^{\star} \mathbf{X} \quad h^{\star} = \mathbf{a}^{\star} \cdot \frac{P_k(\mathbf{x}_{\star})}{\sqrt{d^{\epsilon}}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_{\star} X_{\star}^T}{d} + \kappa_2^2 \frac{H_2(X_{\star}) H_2(X_{\star})^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_{\star}) H_k(X_{\star})^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_{\star}) H_{k+1}(X_{\star})^T}{d^{k+1}} + \dots$$

$$n = O(d^{k\epsilon + \delta})$$

# Spike+ bulk decomposition of Gram matrix

$$\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x} \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_\star X_\star^T}{d} + \kappa_2^2 \frac{H_2(X_\star) H_2(X_\star)^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_\star) H_k(X_\star)^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_\star) H_{k+1}(X_\star)^T}{d^{k+1}} + \dots$$

$$n = O(d^{k\epsilon + \delta})$$

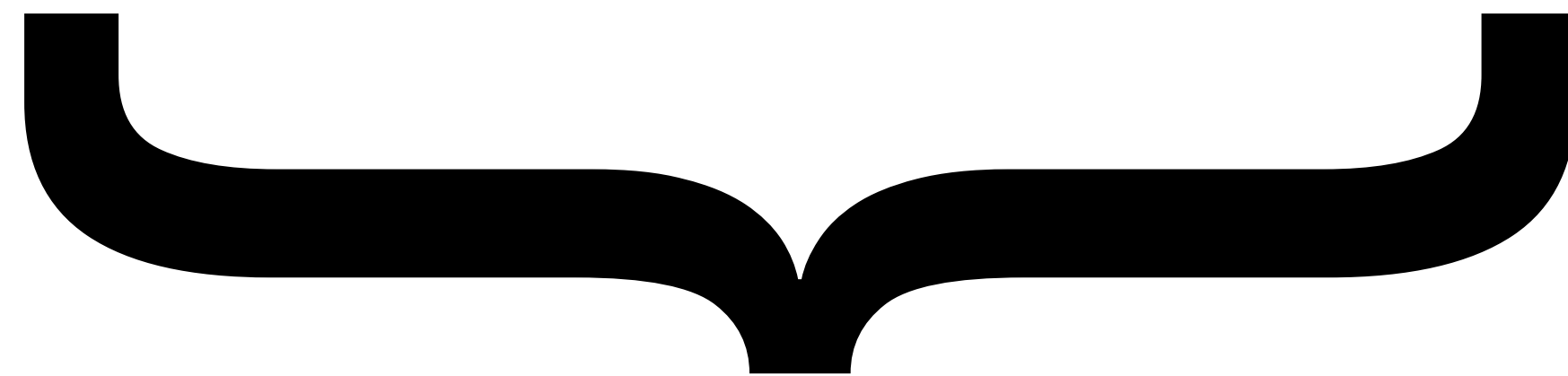
Generalization error of random features and kernel methods:  
hypercontractivity and kernel matrix concentration

Song Mei<sup>\*</sup> Theodor Misiakiewicz<sup>†</sup> Andrea Montanari<sup>†‡</sup>

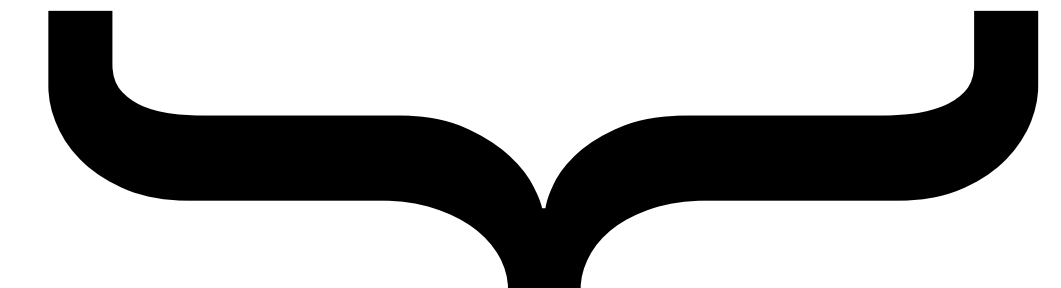
# Spike+ bulk decomposition of Gram matrix

$$\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x} \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_\star X_\star^T}{d} + \kappa_2^2 \frac{H_2(X_\star) H_2(X_\star)^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_\star) H_k(X_\star)^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_\star) H_{k+1}(X_\star)^T}{d^{k+1}} + \dots$$



Concentrates to informative spikes



Noise/ Identity

$$n = O(d^{k\epsilon + \delta})$$

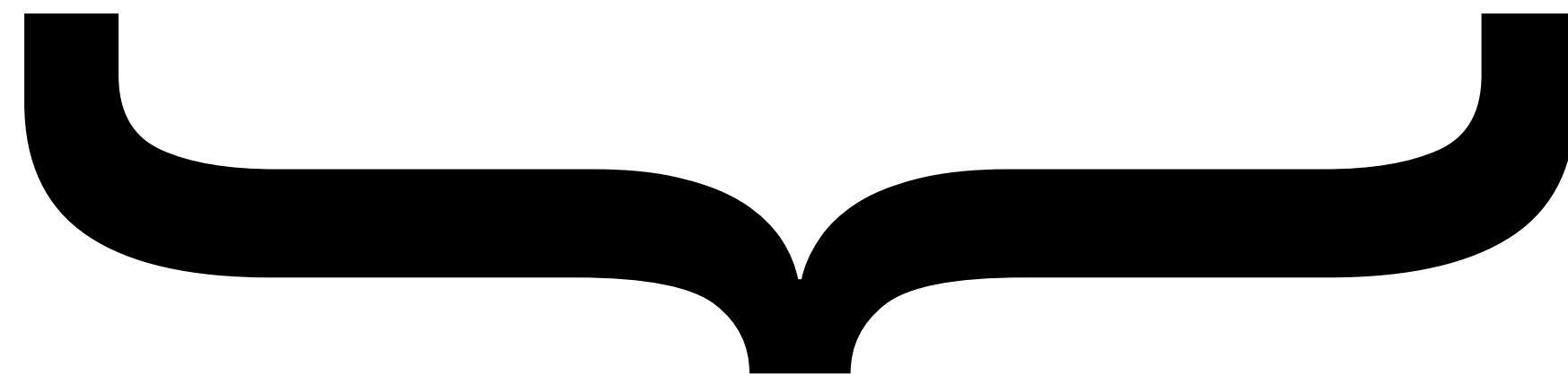
Generalization error of random features and kernel methods:  
hypercontractivity and kernel matrix concentration

Song Mei<sup>\*</sup> Theodor Misiakiewicz<sup>†</sup> Andrea Montanari<sup>†‡</sup>

# Spike+ bulk decomposition of Gram matrix

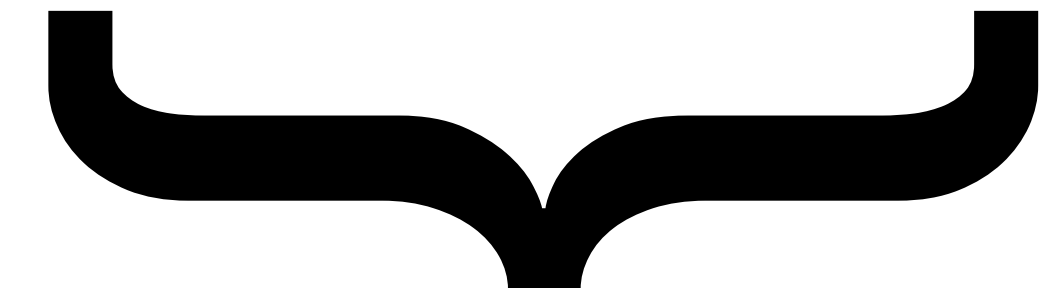
$$\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x} \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_\star X_\star^T}{d} + \kappa_2^2 \frac{H_2(X_\star)H_2(X_\star)^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_\star)H_k(X_\star)^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_\star)H_{k+1}(X_\star)^T}{d^{k+1}} + \dots$$



$$n = O(d^{k\epsilon + \delta})$$

Concentrates to informative spikes



Noise/ Identity

**dominant low-degree term along**  
 $h^\star(\mathbf{x})$

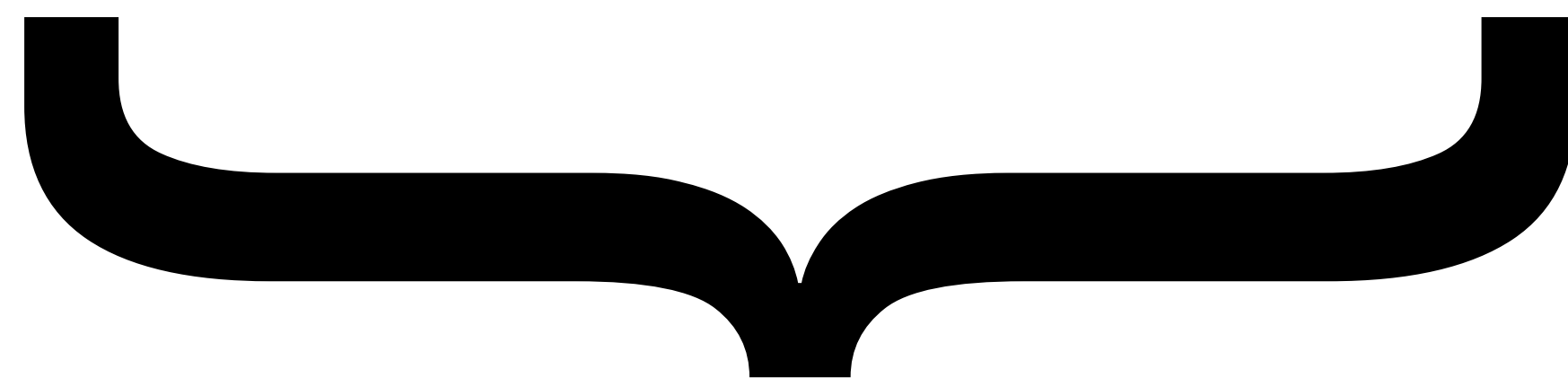
Generalization error of random features and kernel methods:  
hypercontractivity and kernel matrix concentration

Song Mei<sup>\*</sup> Theodor Misiakiewicz<sup>†</sup> Andrea Montanari<sup>†‡</sup>

# Spike+ bulk decomposition of Gram matrix

$$\mathbf{x}_\star = \mathbf{W}^\star \mathbf{x} \quad h^\star = \mathbf{a}^\star \cdot \frac{P_k(\mathbf{x}_\star)}{\sqrt{d^\epsilon}} \in \mathbb{R}$$

$$K \approx \kappa_0^2 \mathbf{1}\mathbf{1}^T + \kappa_1^2 \frac{X_\star X_\star^T}{d} + \kappa_2^2 \frac{H_2(X_\star)H_2(X_\star)^T}{d^2} + \dots + \kappa_k^2 \frac{H_k(X_\star)H_k(X_\star)^T}{d^k} + \kappa_{k+1}^2 \frac{H_{k+1}(X_\star)H_{k+1}(X_\star)^T}{d^{k+1}} + \dots$$



Concentrates to informative spikes



Noise/ Identity

$$n = O(d^{k\epsilon + \delta})$$

**dominant low-degree term along**  
 $h^\star(\mathbf{x})$

Generalization error of random features and kernel methods:  
hypercontractivity and kernel matrix concentration

Song Mei<sup>\*</sup> Theodor Misiakiewicz<sup>†</sup> Andrea Montanari<sup>†‡</sup>

**Caveat:**

With Gaussian inputs, some radial degree  $k$  polynomials require less than  $O(d^k)$  samples.

# Fitting the last layer



# Fitting the last layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(W_1 \mathbf{x}) \approx c \mathbf{w}_3 \mathbf{h}^\star(\mathbf{x})$$

# Fitting the last layer

$$\mathbf{h}(\mathbf{x}) = W_2 \sigma(\mathbf{W}_1 \mathbf{x}) \approx c \mathbf{w}_3 \mathbf{h}^\star(\mathbf{x})$$

- Reduction to Kernel on low-dimensional features.
- $K(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(c \mathbf{w} \mathbf{h}^\star(\mathbf{x}_1) + b) \sigma(c \mathbf{w} \mathbf{h}^\star(\mathbf{x}_2) + b)]$ .
- $\hat{W}_3 \approx \text{KRR}(K(\cdot, \cdot), \mathbf{X}, \mathbf{y})$ .
- Sample complexity/width now dimension-independent.

# Takeaways

# Takeaways

**Hierarchical functions with robustness w.r.t intermediate features  
allow exploitation of depth through dimension reduction**

# Takeaways

**Hierarchical functions with robustness w.r.t intermediate features  
allow exploitation of depth through dimension reduction**

Do we need narrowing of networks? No, consider:

$$f^\star(\mathbf{x}) = g^\star \left( \frac{\mathbf{a}_1^{\star\top} P_2(W_1^\star \mathbf{x})}{\sqrt{d}}, \dots, \frac{\mathbf{a}_m^{\star\top} P_2(W_m^\star \mathbf{x})}{\sqrt{d}} \right), \mathbf{x} \in \mathbb{R}^d, m = \mathcal{O}(d)$$

# Takeaways

**Hierarchical functions with robustness w.r.t intermediate features  
allow exploitation of depth through dimension reduction**

Do we need narrowing of networks? No, consider:

$$f^\star(\mathbf{x}) = g^\star \left( \underbrace{\frac{\mathbf{a}_1^{\star\top} P_2(W_1^\star \mathbf{x})}{\sqrt{d}}, \dots, \frac{\mathbf{a}_m^{\star\top} P_2(W_m^\star \mathbf{x})}{\sqrt{d}}}_{\text{“approximately independent” } \mathcal{O}(d) \text{ features in } \mathcal{O}(d^2) \text{ space}} \right), \mathbf{x} \in \mathbb{R}^d, m = \mathcal{O}(d)$$

“approximately independent”  $\mathcal{O}(d)$  features in  $\mathcal{O}(d^2)$  space



# Thanks to my collaborators!

