

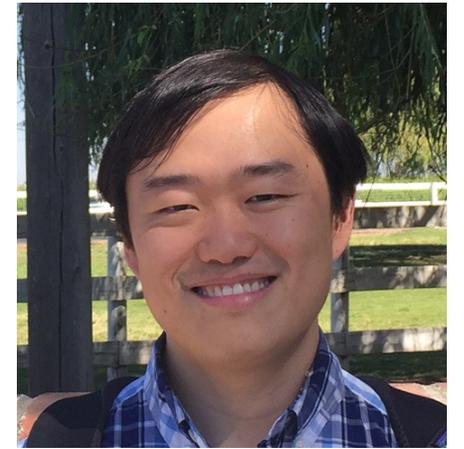
# Understanding Deep Models With Teacher-student Setting.

Yuandong Tian

Research Scientist and Manager  
Facebook AI Research



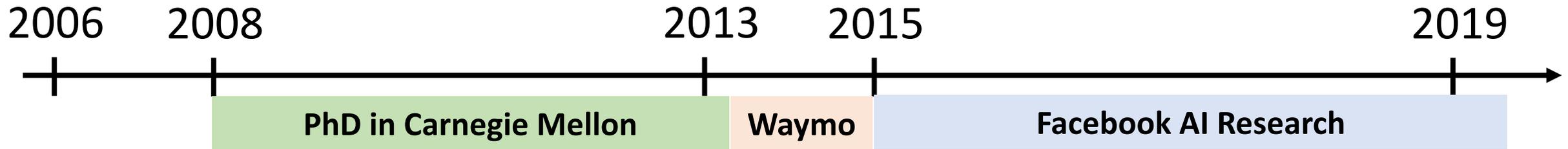
# Career Path



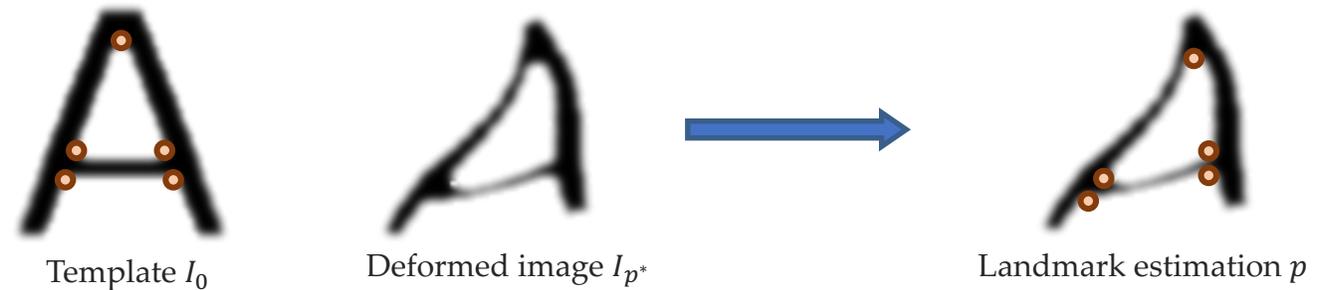
**Theoretical Understanding of Models and Algorithms**

**Computer Vision**

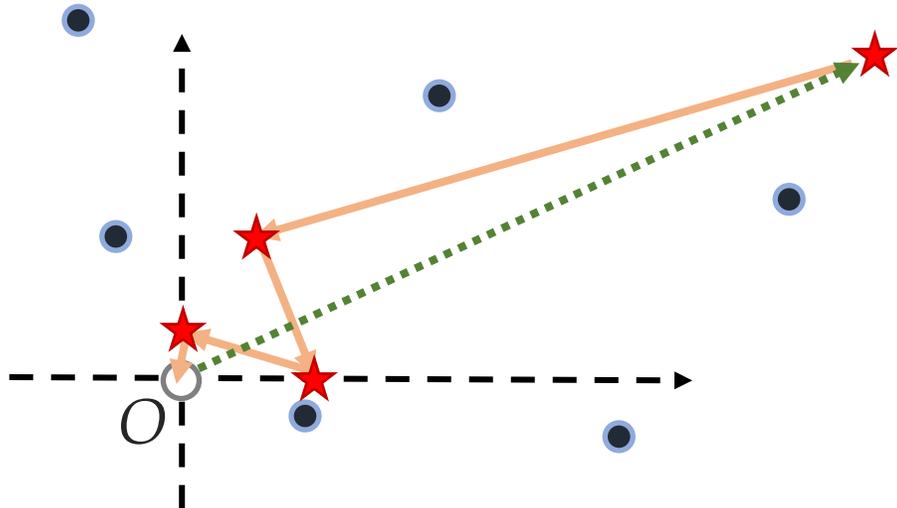
**Reinforcement Learning**



# Data-Driven Descent (PhD work)



To achieve  $\|p - p^*\| \leq \epsilon$



Method	Sample complexity
Gradient descent	Local optimality
Nearest Neighbor	$O(\epsilon^{-d})$
[Y. Tian and S. Narasimhan, CVPR 10]	$O(C^d \log \epsilon^{-1})$
[Y. Tian and S. Narasimhan, ICCV 13, <i>Marr Prize Honorable Mention</i> ]	$O(C_1^d + C_2 \log \epsilon^{-1})$

# DarkForestGo



DarkForest versus Koichi Kobayashi (9p)

*[Better Computer Go Player with Neural Network and Long-term Prediction, Y. Tian and Y. Zhu, ICLR 2016]*

# How Facebook's AI Researchers Built a Game-Changing Go Engine

The best human players easily beat the best computer-based Go engines. That looks set to change thanks to a new approach pioneered by Facebook's artificial intelligence researchers.

by **Emerging Technology** from the arXiv

Dec 4, 2015

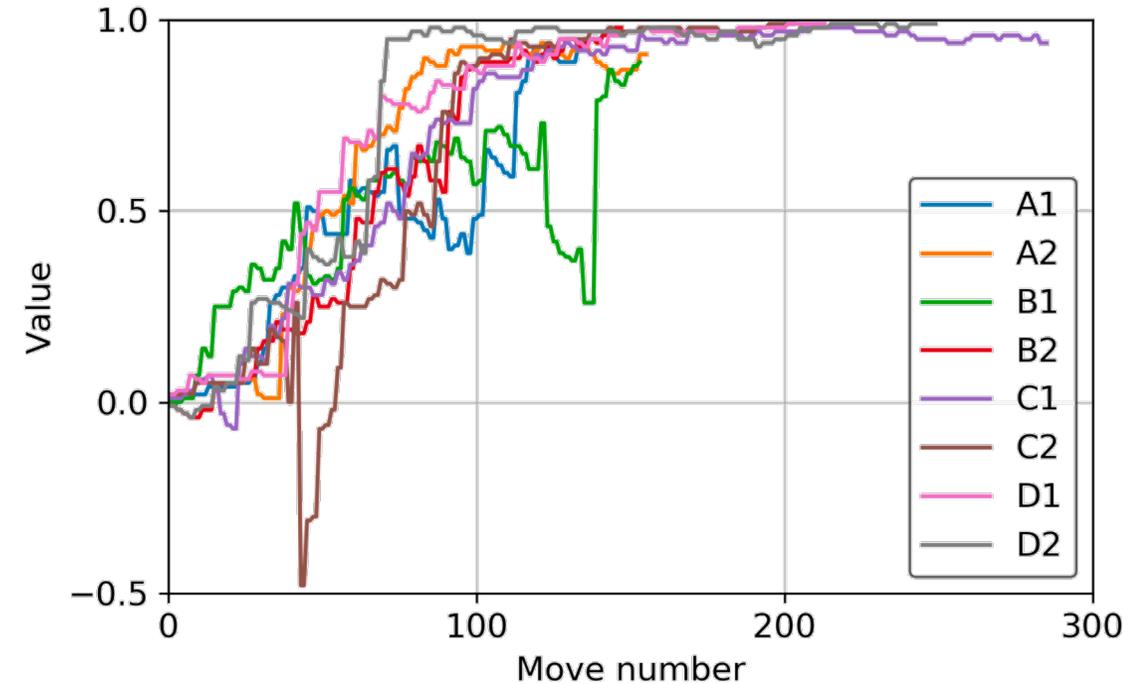
**One of the last bastions of human mastery over computers is the game of Go** —the best human players beat the best Go engines with ease.

# ELF OpenGo

## Vs top professional players

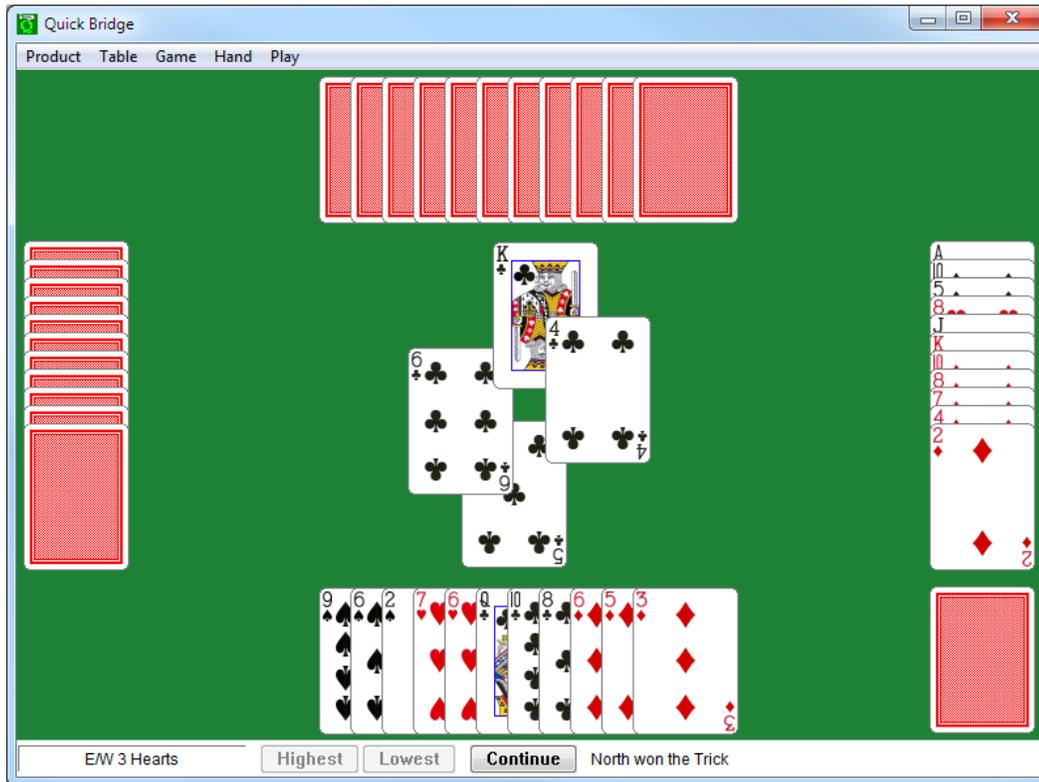
Name (rank)	ELO (world rank)	Result
Kim Ji-seok	3590 (#3)	5-0
Shin Jin-seo	3570 (#5)	5-0
Park Yeonghun	3481 (#23)	5-0
Choi Cheolhan	3466 (#30)	5-0

Single GPU, 80k rollouts, 50 seconds  
Offer unlimited thinking time for the players



*[ELF OpenGo: An Analysis and Open Reimplementation of AlphaZero, Y. Tian et al, ICML 2019]*

# Contract Bridge



SoTA performance on Bridge Bidding

A new theoretical framework for multi-agent collaborative games



Yuandong Tian

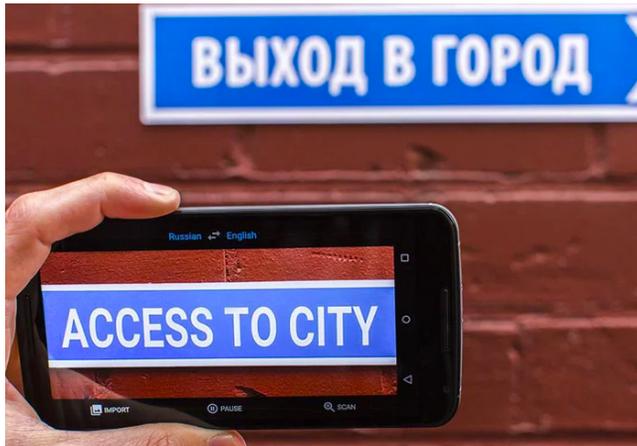
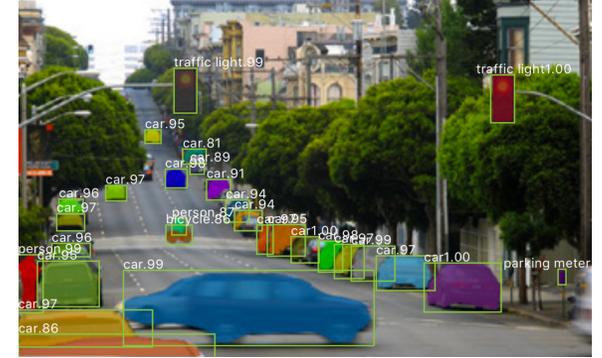


Qucheng Gong



Tina Jiang

# Great Empirical Success



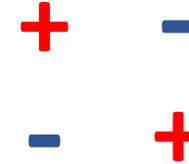
# How do deep models work?



# Three Major Problems

Understanding how  
Deep Models work

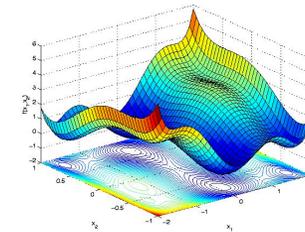
## Expressibility



“Neural Network is a universal approximator”

“Deep Models can express functions more efficiently than shallow ones”

## Optimization

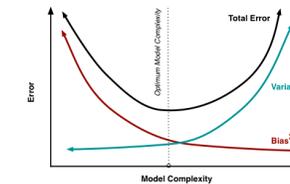


“Gradient vanishing/exploding”

“Gradient Descent might get stuck at saddle point / local minima”

“Can GD/SGD go to global optima? How fast?”

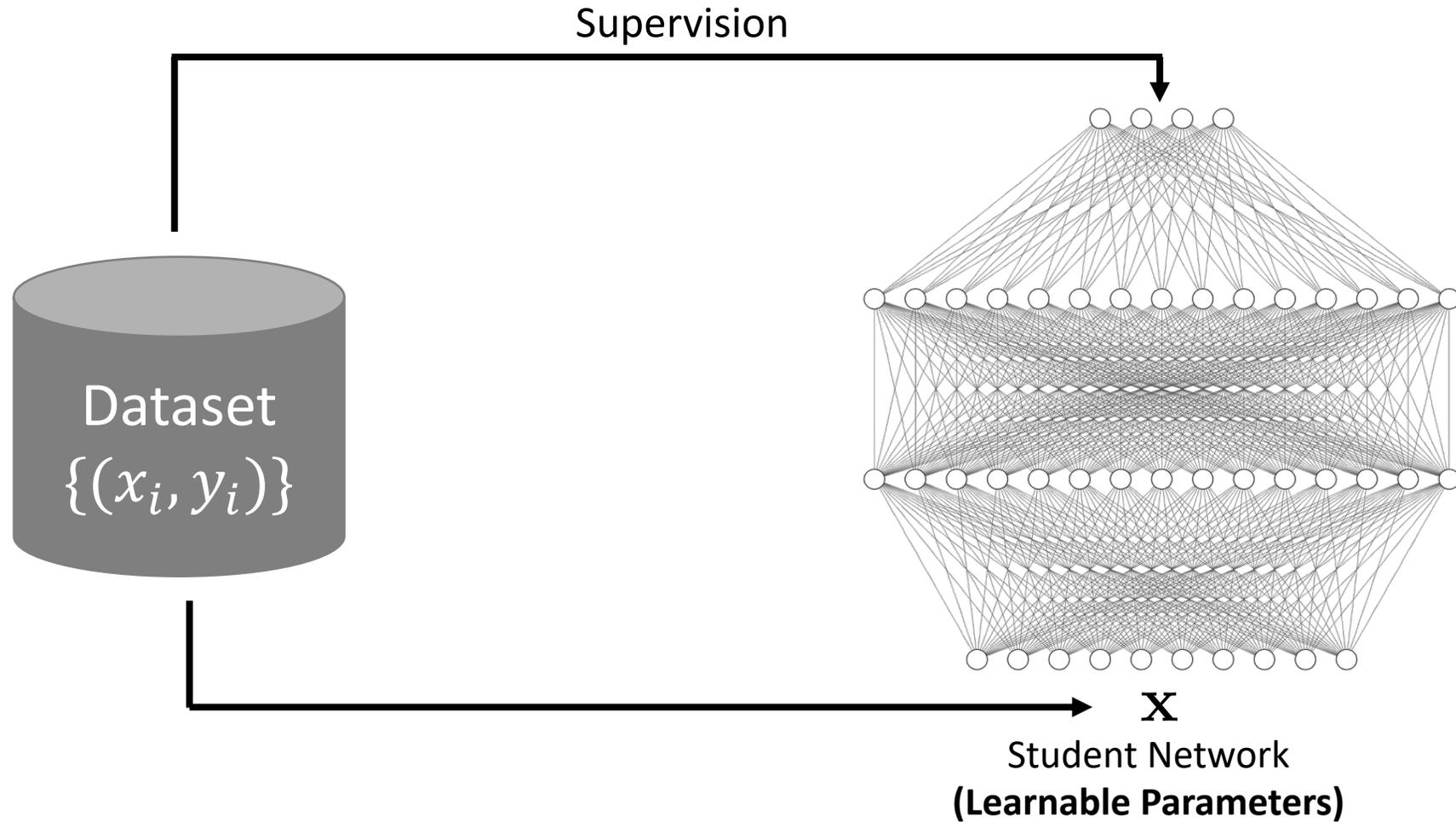
## Generalization



“Does zero training error often lead to overfitting?”

“More parameters might lead to overfitting.”

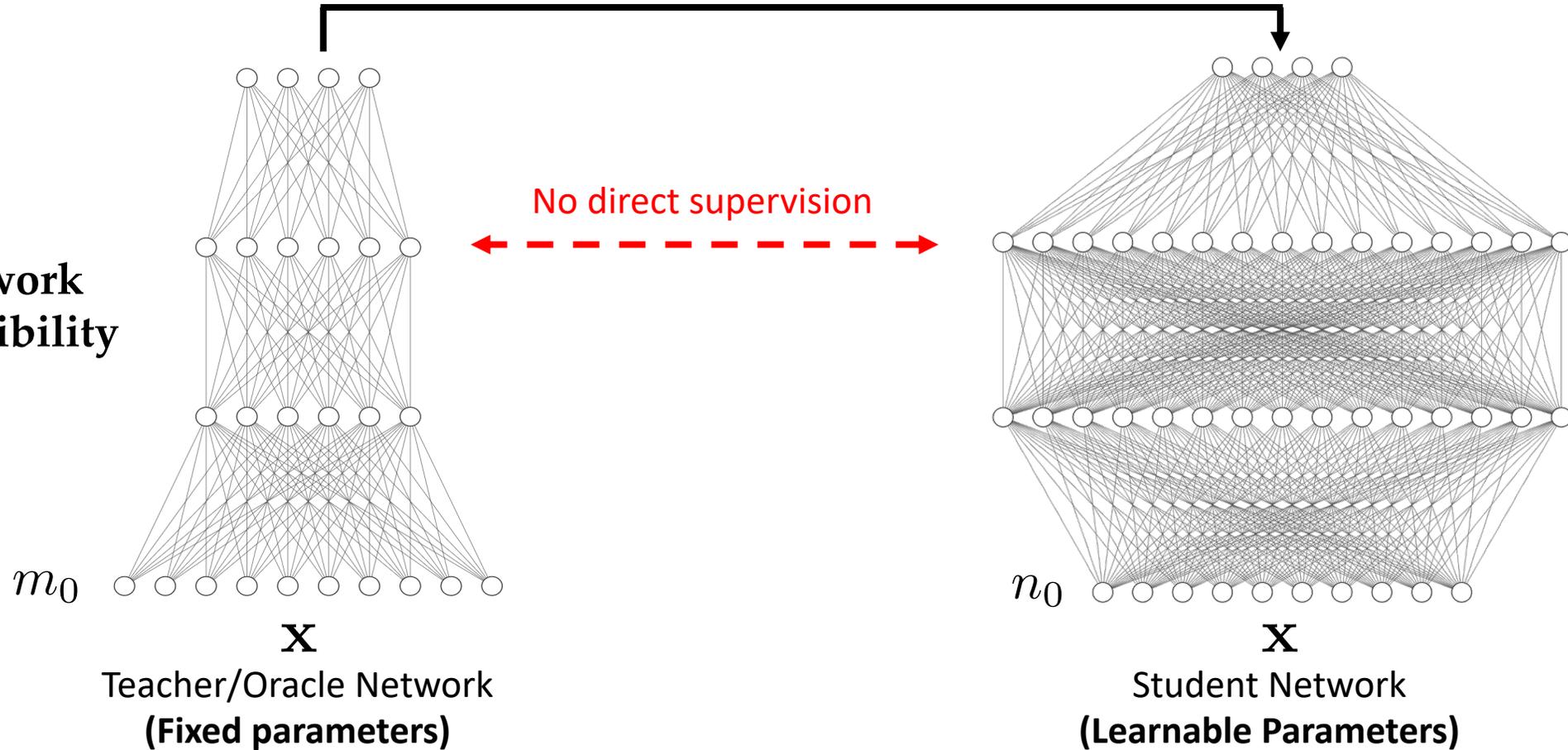
# Supervised Learning



# Student-Teacher Setting

Supervision

By Network  
Expressibility



# Why Student-Teacher Setting?

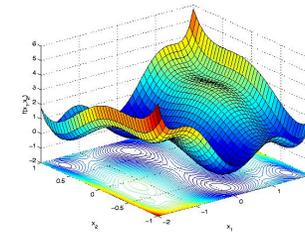
Understanding how  
Deep Models work

**Expressibility**

**+** **-**  
**-** **+**

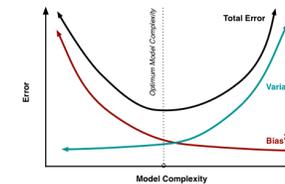
Provide a target function with bounded complexity

**Optimization**



Study fine dynamics behaviors by comparing with teacher

**Generalization**



**Our focus** → **Student Specialization** yields generalization

# Old History of Teacher-Student Setting

$$\epsilon(\mathbf{J}) = \frac{1}{2} \langle |f(\mathbf{J}, \xi) - f(\mathbf{B}, \xi)|^2 \rangle_{\xi} \quad f(\mathbf{J}, \xi) = \sum_{i=1}^K \sigma(\mathbf{J}_i \cdot \xi)$$

Study when the input dimension  $n_0 = m_0 \rightarrow +\infty$  (i.e., **thermodynamics limits**)

In some situations, student nodes are “**specialized**” to teacher node

One layer of trainable parameters

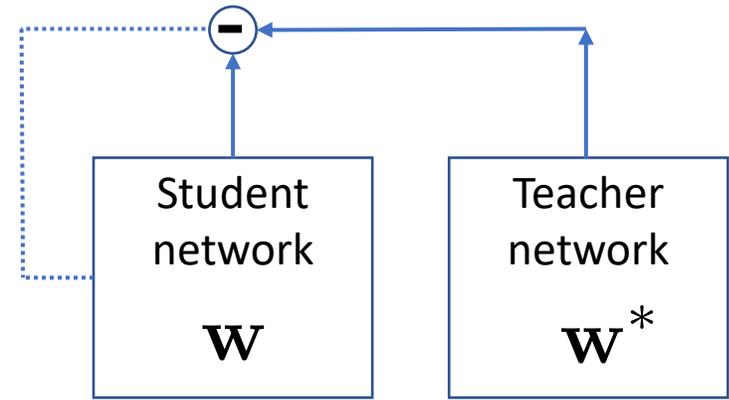
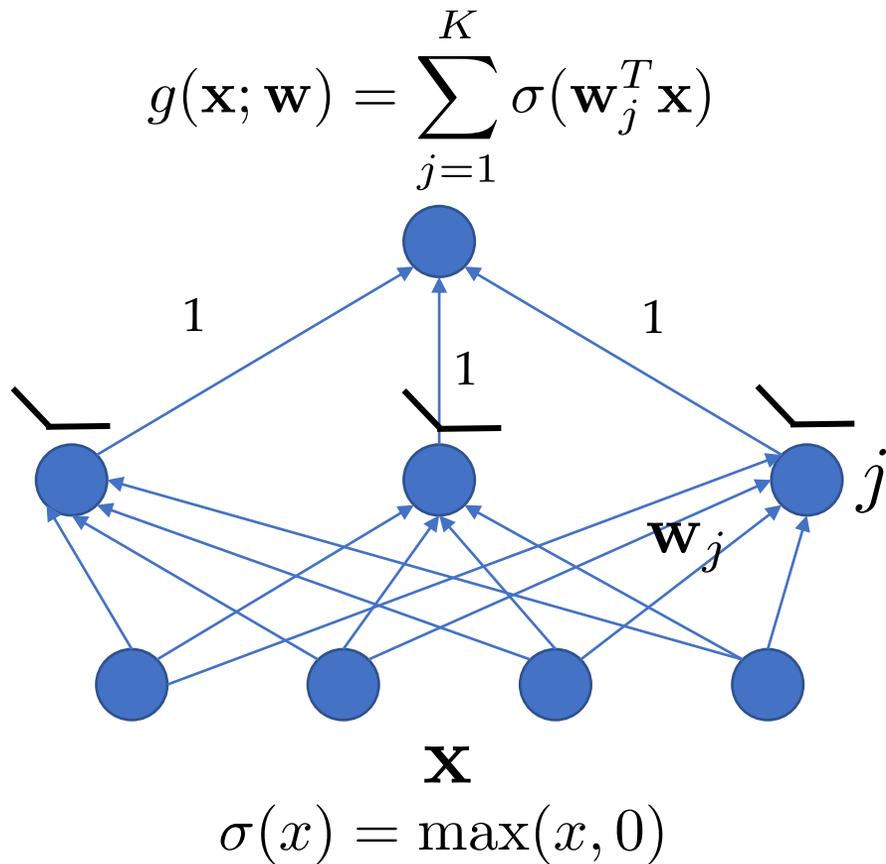
Nonlinear function  $\sigma(x) = \text{erf}(x / 2)$

Locally linearized analysis around symmetry breaking plane and final solution

*[On-line learning in soft committee machines, Saad & Solla, Phys. Rev 1995]*

*[S. Goldt et al, Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup, NeurIPS 2019]*

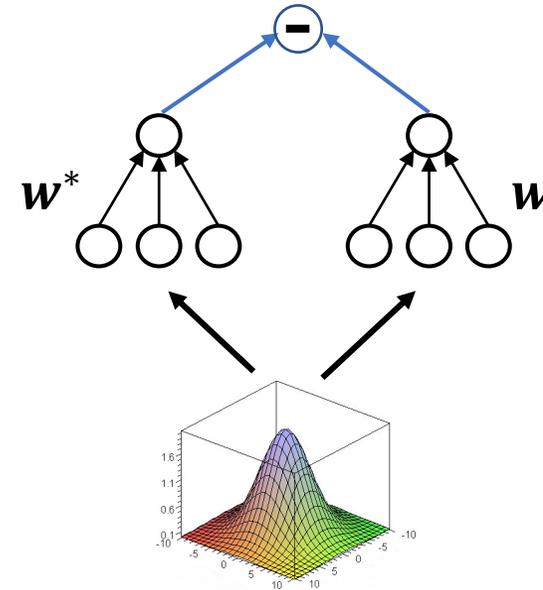
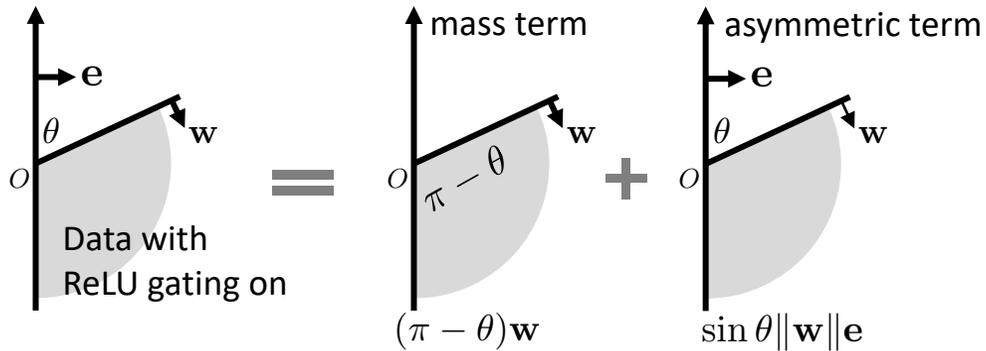
# Simplest Teacher Student Setting: ReLU networks with Gaussian Inputs



$$J(\mathbf{w}) = \frac{1}{2} \|g(X; \mathbf{w}^*) - g(X; \mathbf{w})\|^2$$

We focus on *population gradient*

# An Analytic Formula



Close-form Population Gradient:

$$\mathbb{E} [\nabla_{\mathbf{w}} J] = \frac{N}{2} (\mathbf{w} - \mathbf{w}^*) + \frac{N}{2\pi} \left( \theta \mathbf{w}^* - \frac{\|\mathbf{w}_{\mathbf{X}}^*\|}{\|\mathbf{w}\|} \sin \theta \mathbf{w} \right)$$

Linear component.

Global convergence if this is the only term

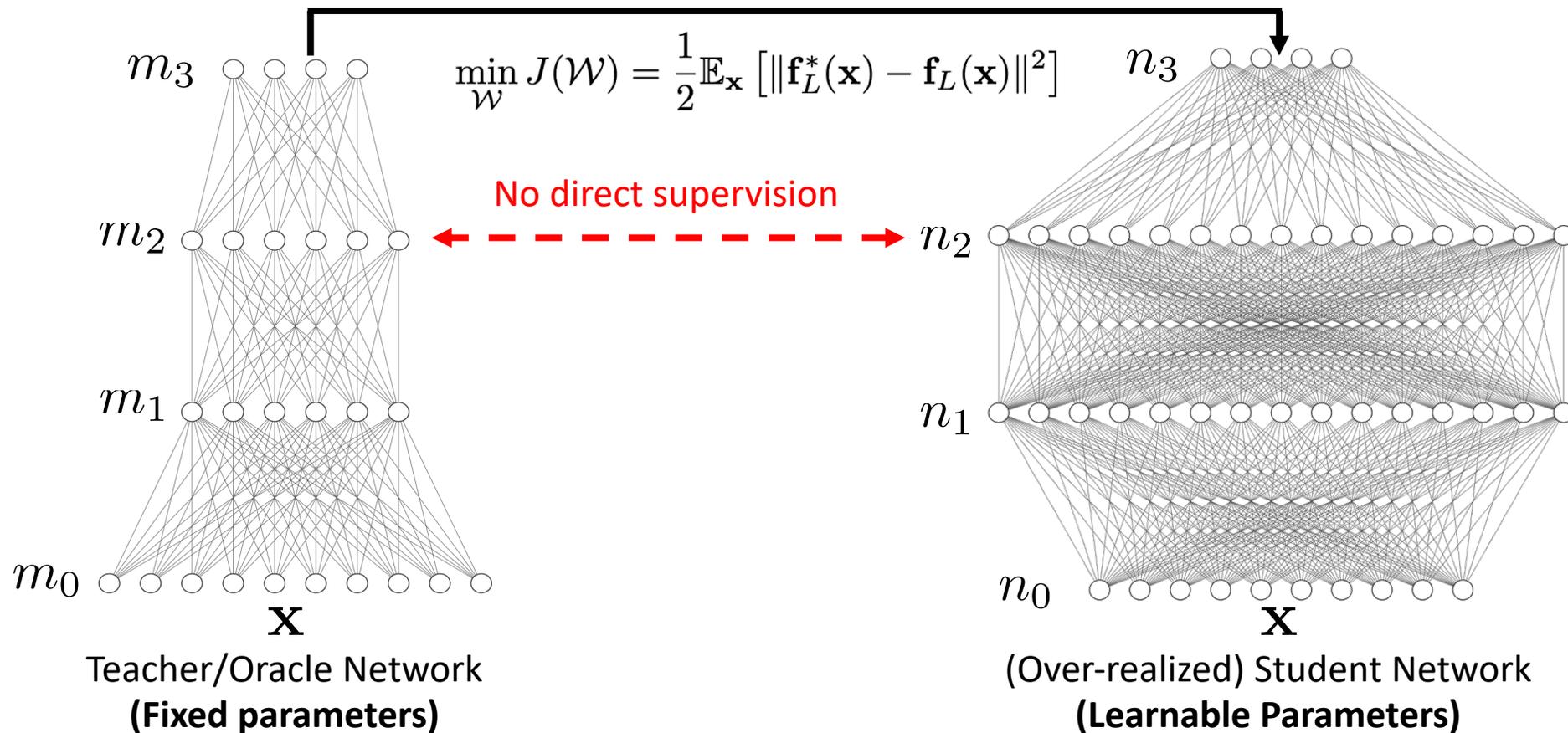
Nonlinear component

due to ReLU gating

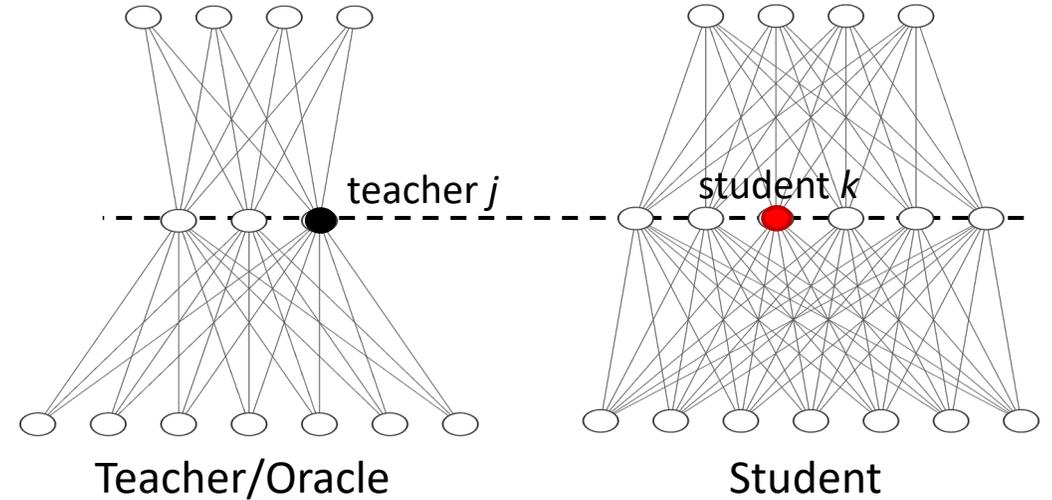
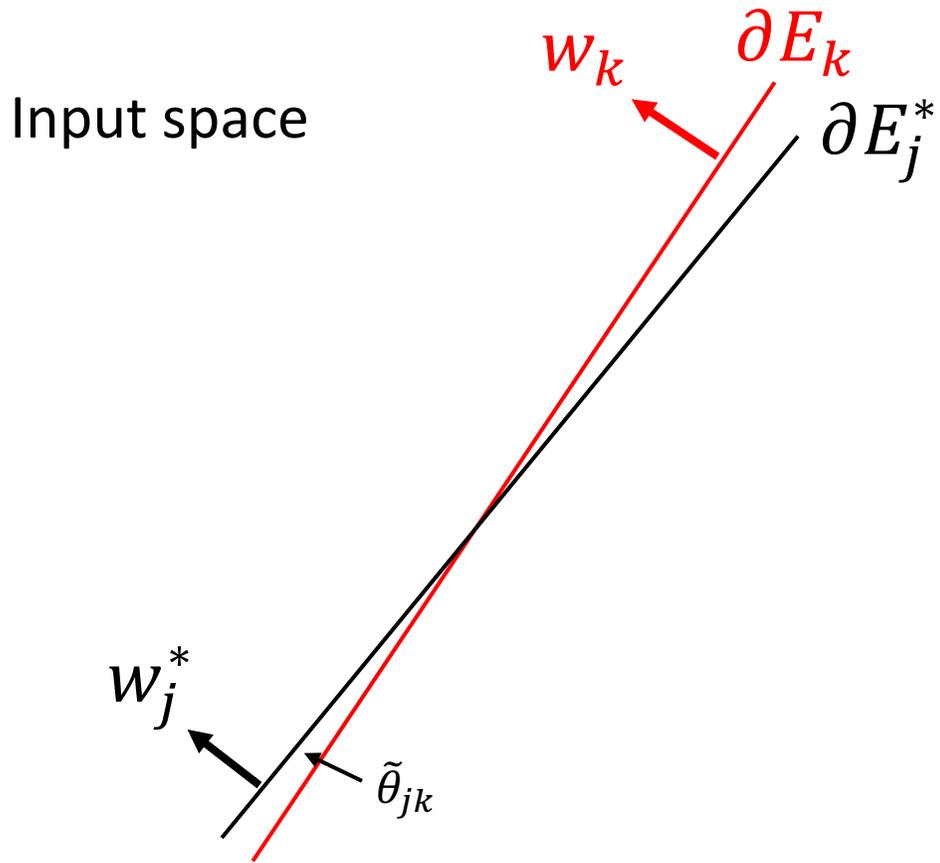
# Multi-layer ReLU network

1. Finite  $m_0$  and  $n_0$
2. Works for  $n_i \geq m_i$   
(no crazy overparameterization)

*Different From Neural Tangent Kernel*



# Student Specialization



$\partial E_k$ : Boundary of node  $k$

$\partial E_j^*$ : Boundary of teacher node  $j$

**$\epsilon$ -alignment:**  $\sin \tilde{\theta}_{jk} \leq \epsilon$  and  $|b_j - b_k^*| \leq \epsilon$

# Main Question

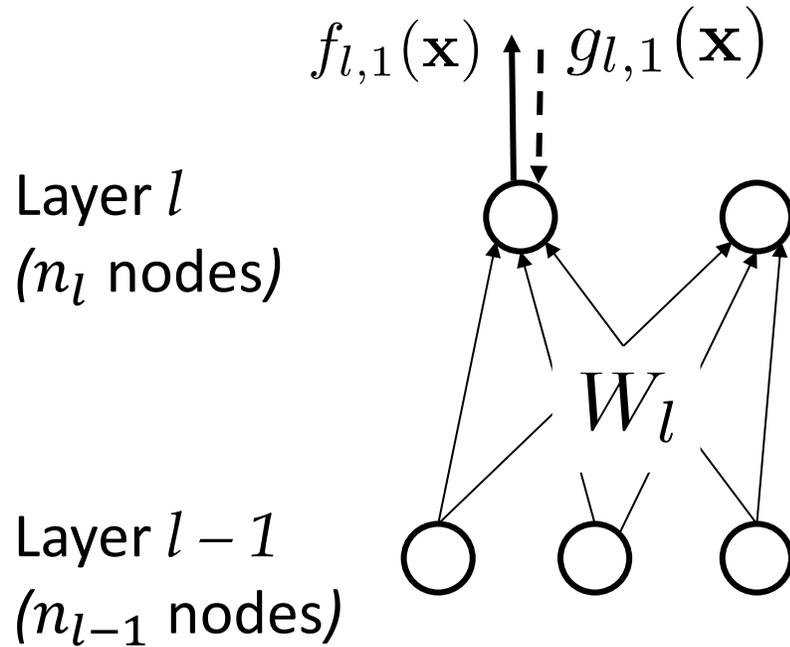
Small gradient  
at every training sample  
during training



Student aligns  
with the teacher

→ Small training error leads to good generalization

# Notation



**Activation**

$$\mathbf{f}_l(\mathbf{x}) = \begin{bmatrix} f_{l,1}(\mathbf{x}) \\ f_{l,2}(\mathbf{x}) \end{bmatrix}$$

**Gradient**

$$\mathbf{g}_l(\mathbf{x}) = \begin{bmatrix} g_{l,1}(\mathbf{x}) \\ g_{l,2}(\mathbf{x}) \end{bmatrix}$$

**Weight update rule:**  $\dot{W}_l = \mathbb{E}_{\mathbf{x}} [\mathbf{f}_{l-1}(\mathbf{x}) \mathbf{g}_l^\top(\mathbf{x})]$

GD: expectation taken over the entire dataset

SGD: expectation taken over a batch

# Lemma1: Recursive Gradient Rule

For layer  $l$ , there exists  $A_l(x)$  and  $B_l(x)$  so that:

$$\mathbf{g}_l(\mathbf{x}) = D_l(\mathbf{x}) [A_l(\mathbf{x})\mathbf{f}_l^*(\mathbf{x}) - B_l(\mathbf{x})\mathbf{f}_l(\mathbf{x})]$$

Student gradient                                            Teacher mixture                      Student mixture

Student gating

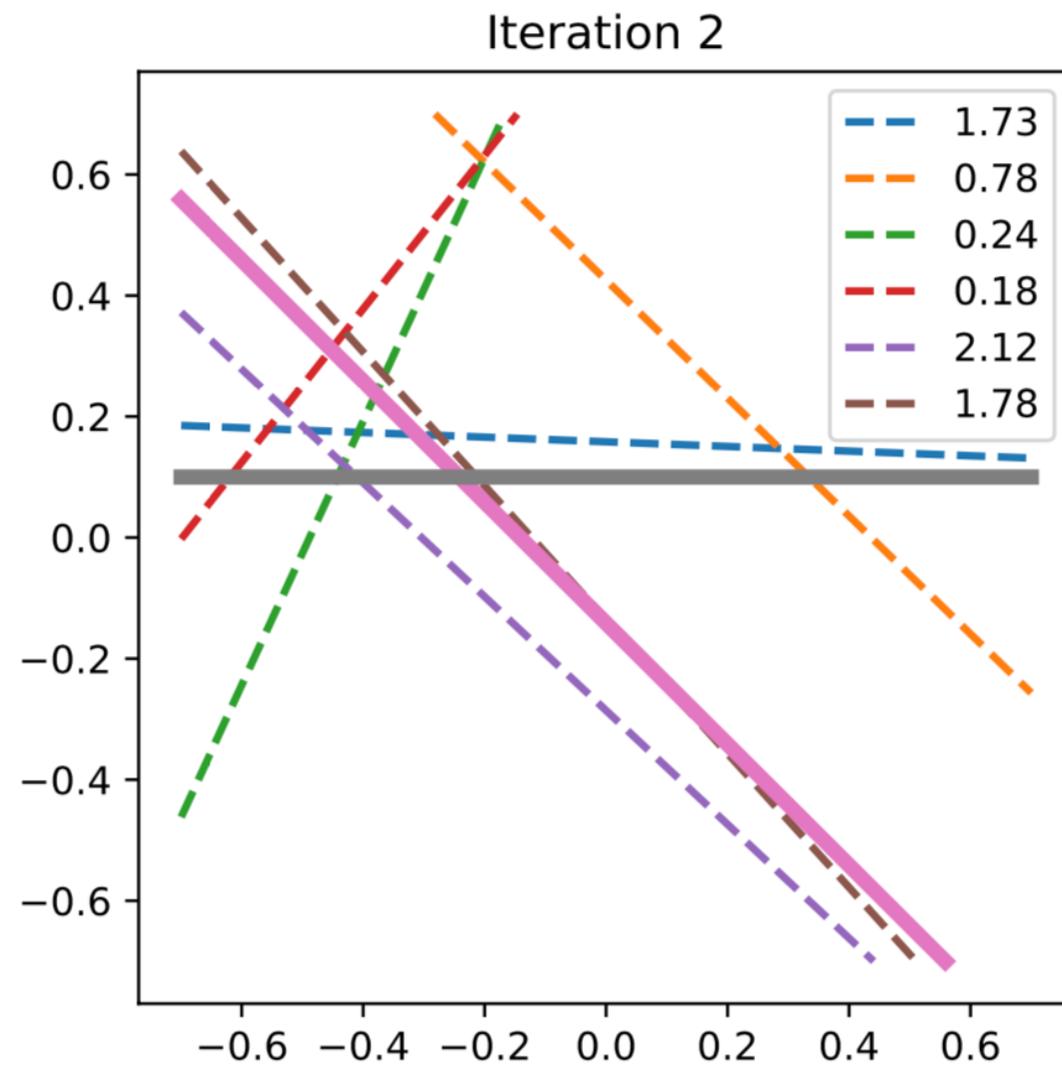
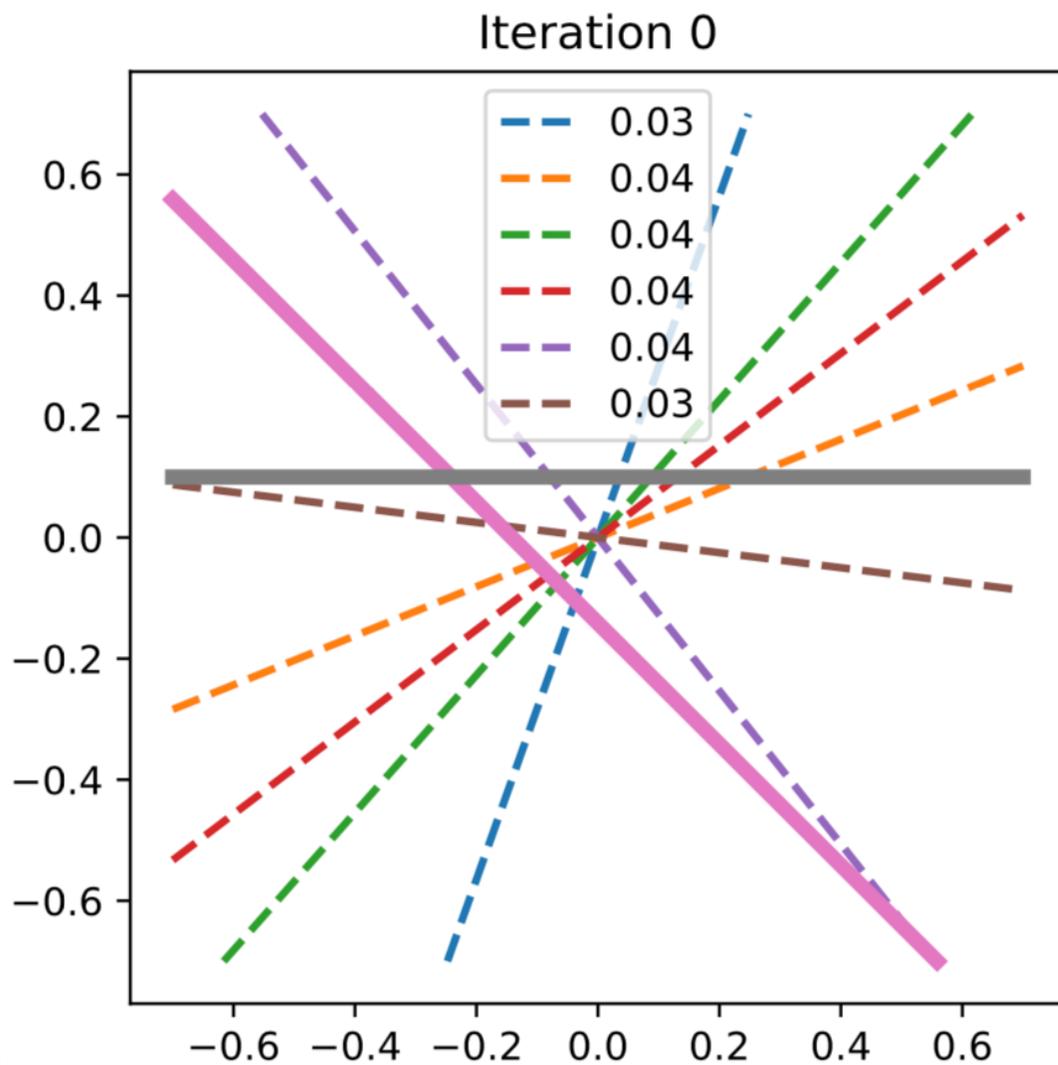
$A_l(x)$  and  $B_l(x)$  are **piece-wise constant**.

A Demonstrative Case:

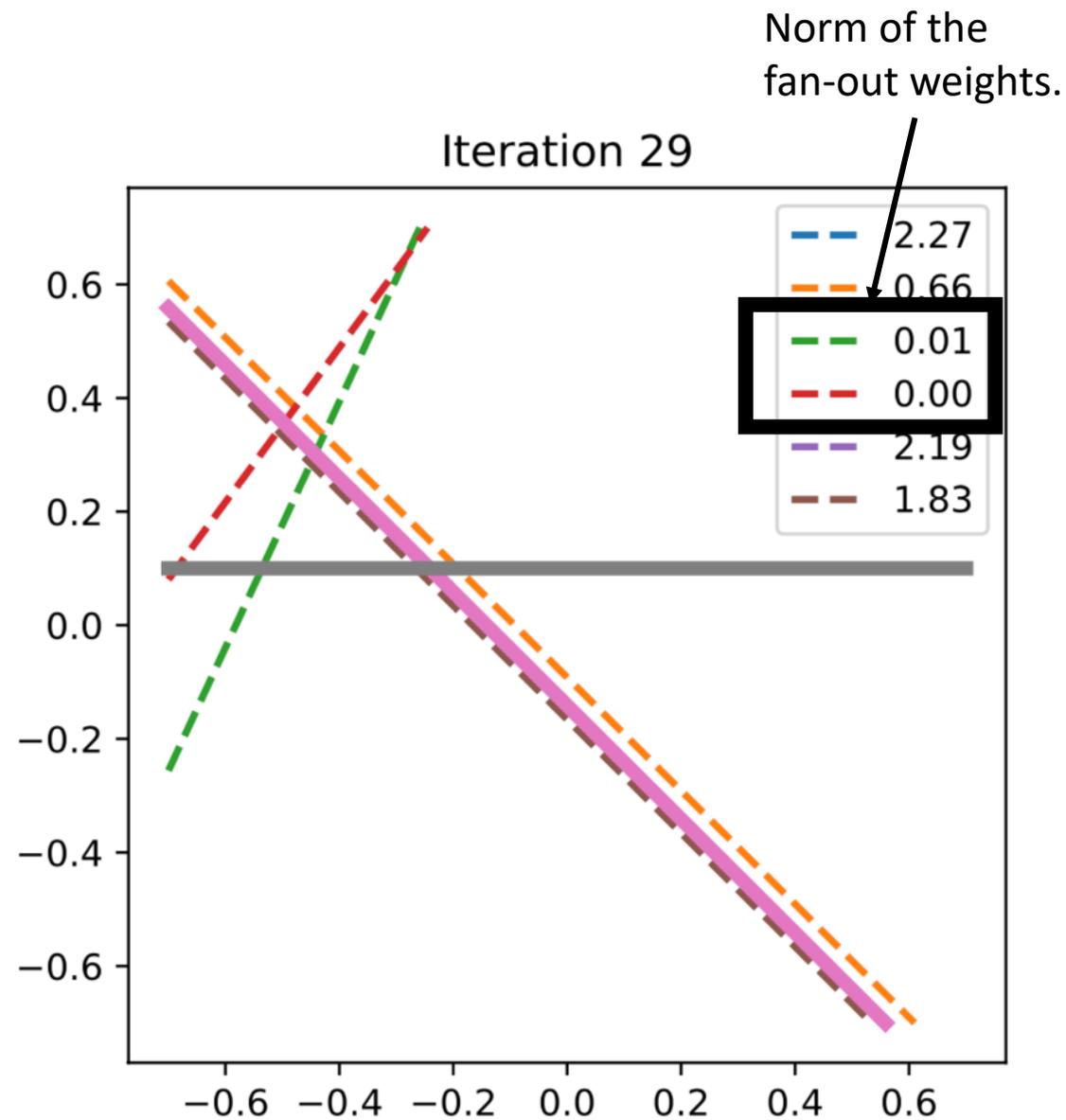
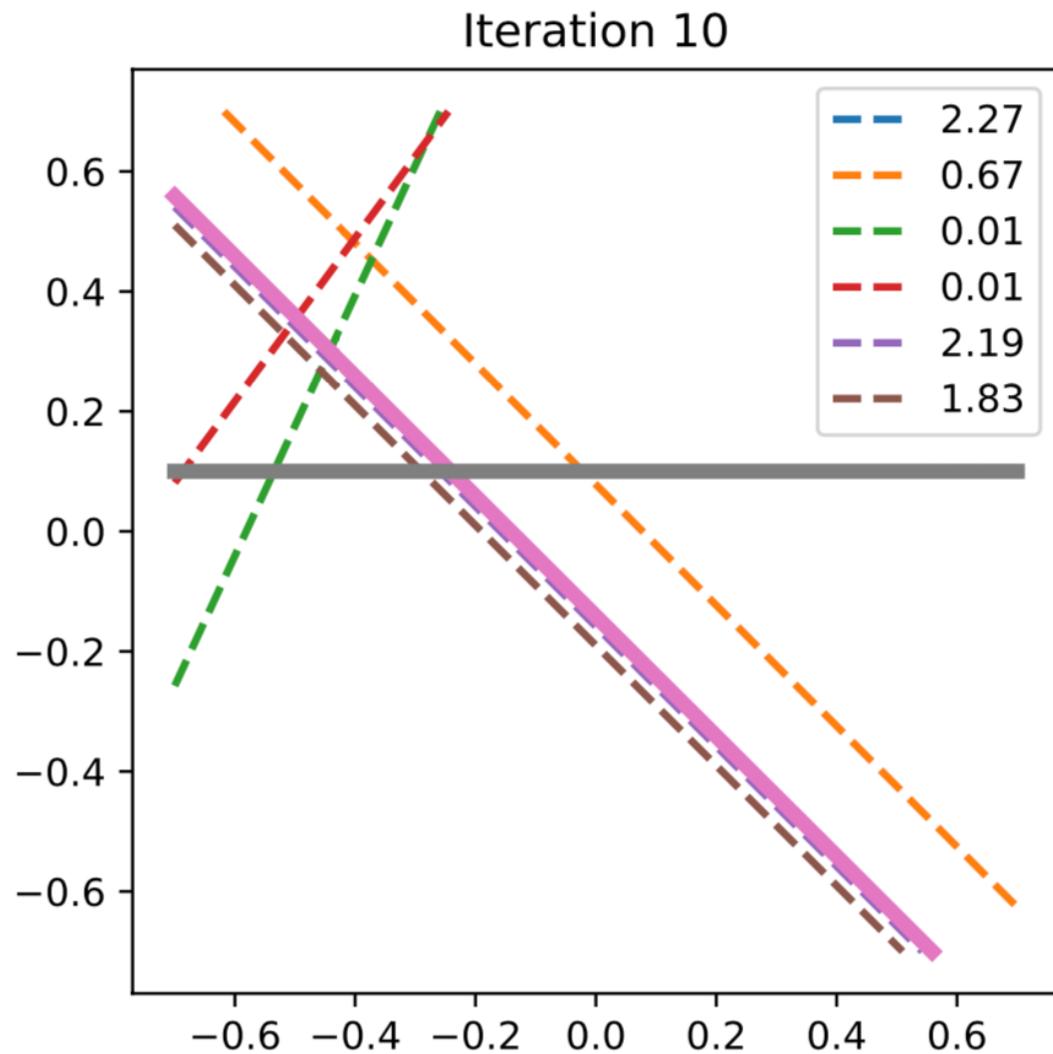
**Two-layer** Network, **Zero Gradient** and  
**Infinite** Samples

# Simple 2D experiments

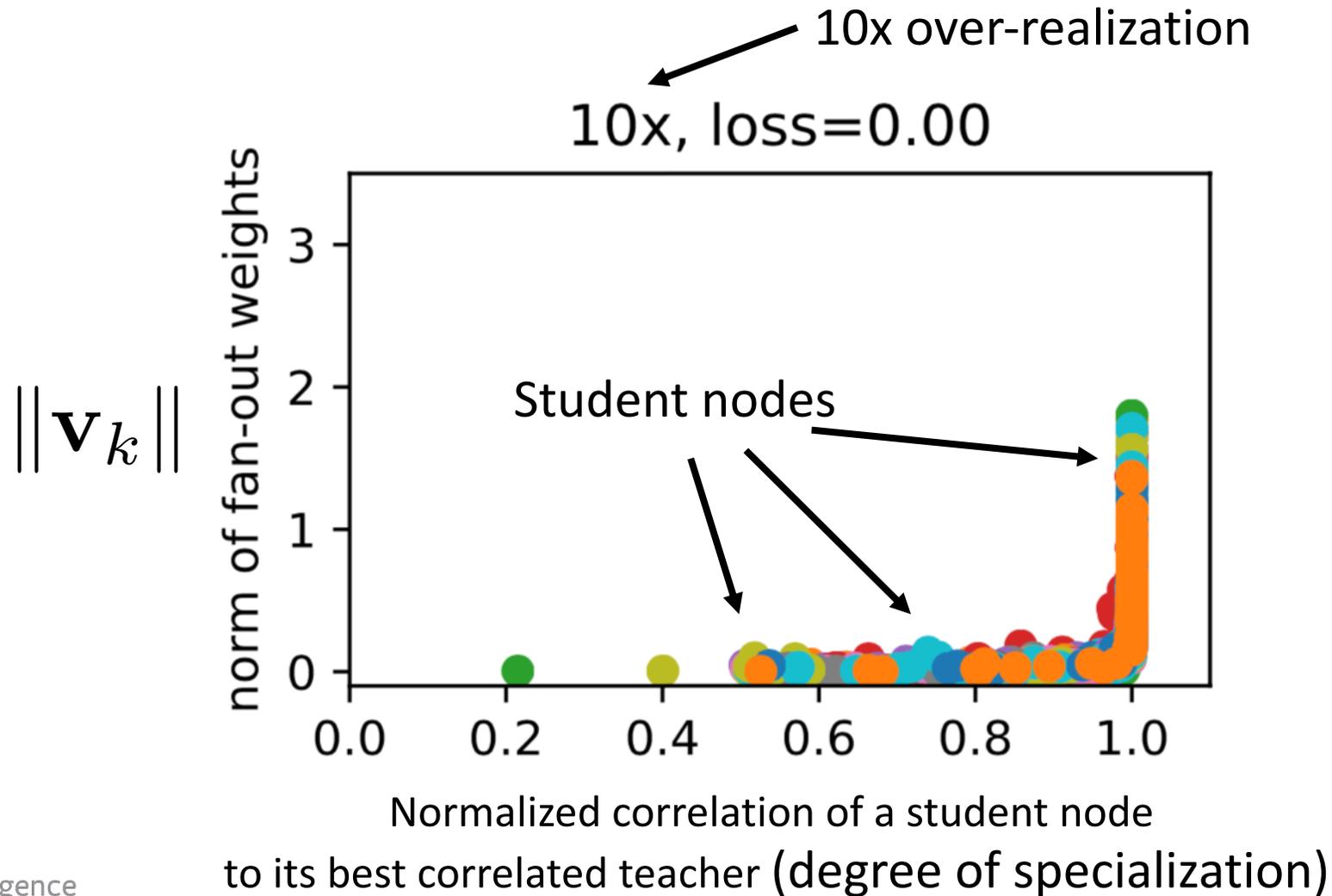
Student Boundary  
Teacher Boundary



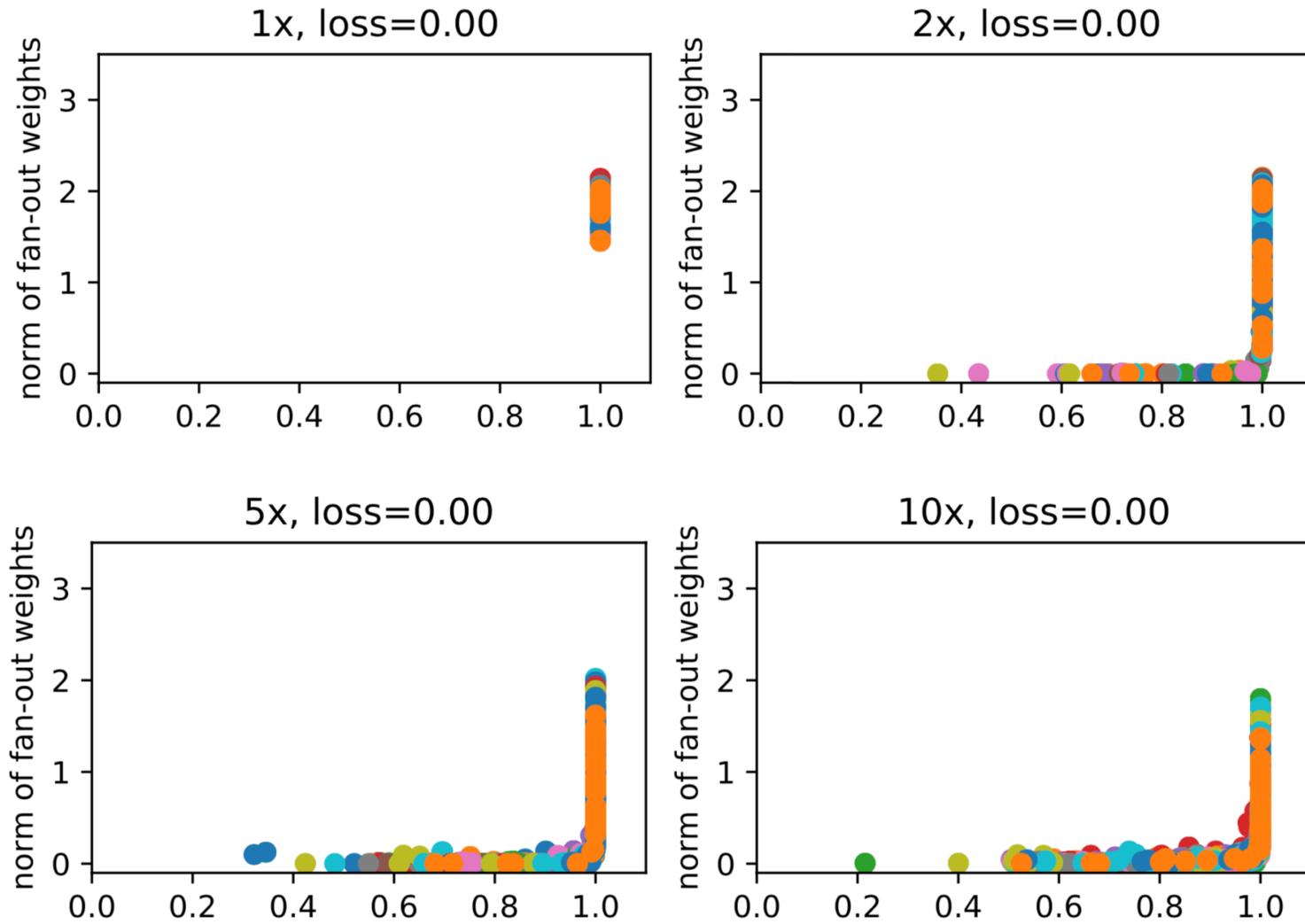
# Simple 2D experiments



# L-shape curve at convergence

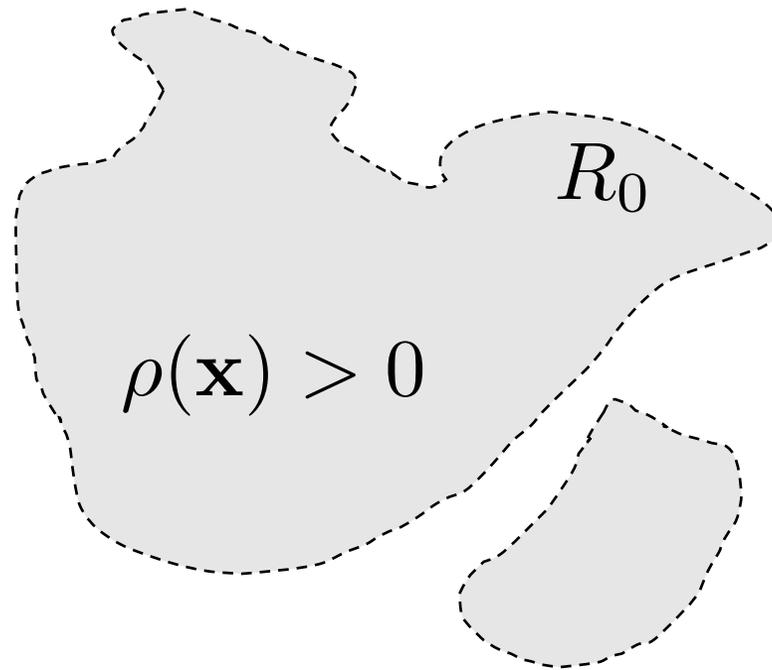


# L-shape curve at convergence



# Assumption of the dataset

No parametrized assumptions



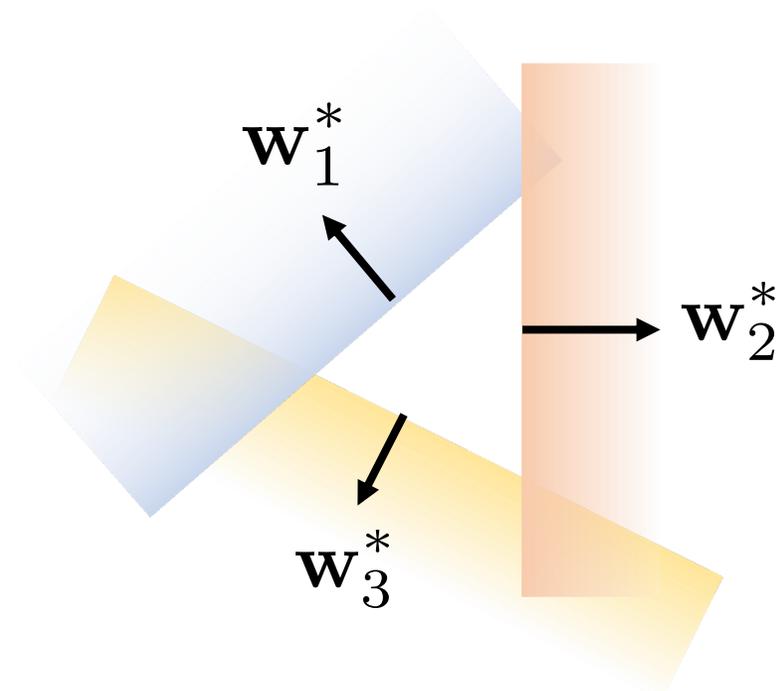
Infinite dataset!

(Region needs to have interiors)

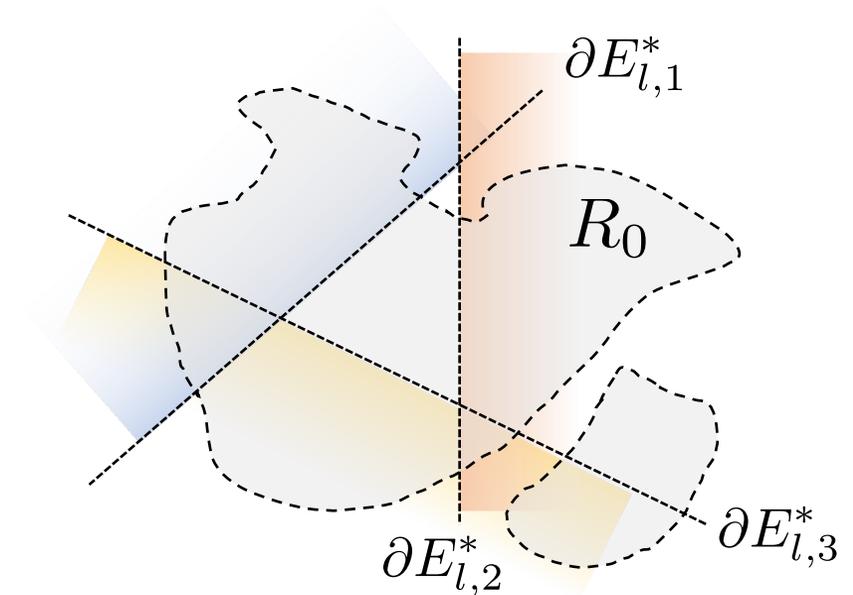
***Full rank***

# Assumptions on Teacher Network

- Cannot reconstruct arbitrary teachers
  - e.g., all ReLU nodes are dead

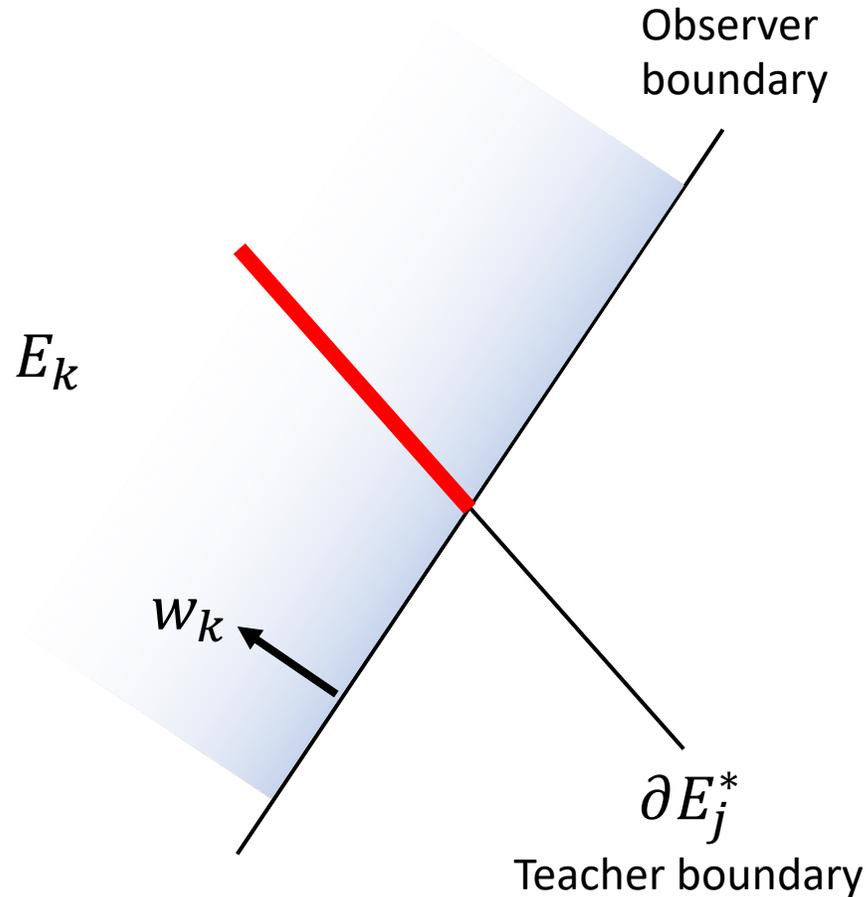


Distinct teacher nodes



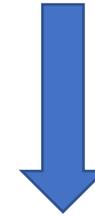
Teacher's ReLU boundary are **visible** in the dataset

# Definition of “Observation”



$E_k$ : Activation region of node  $k$

$$\partial E_j^* \cap E_k \neq \emptyset$$



Teacher  $j$  is **observed** by a student  $k$

# Main results: Alignment could happen!

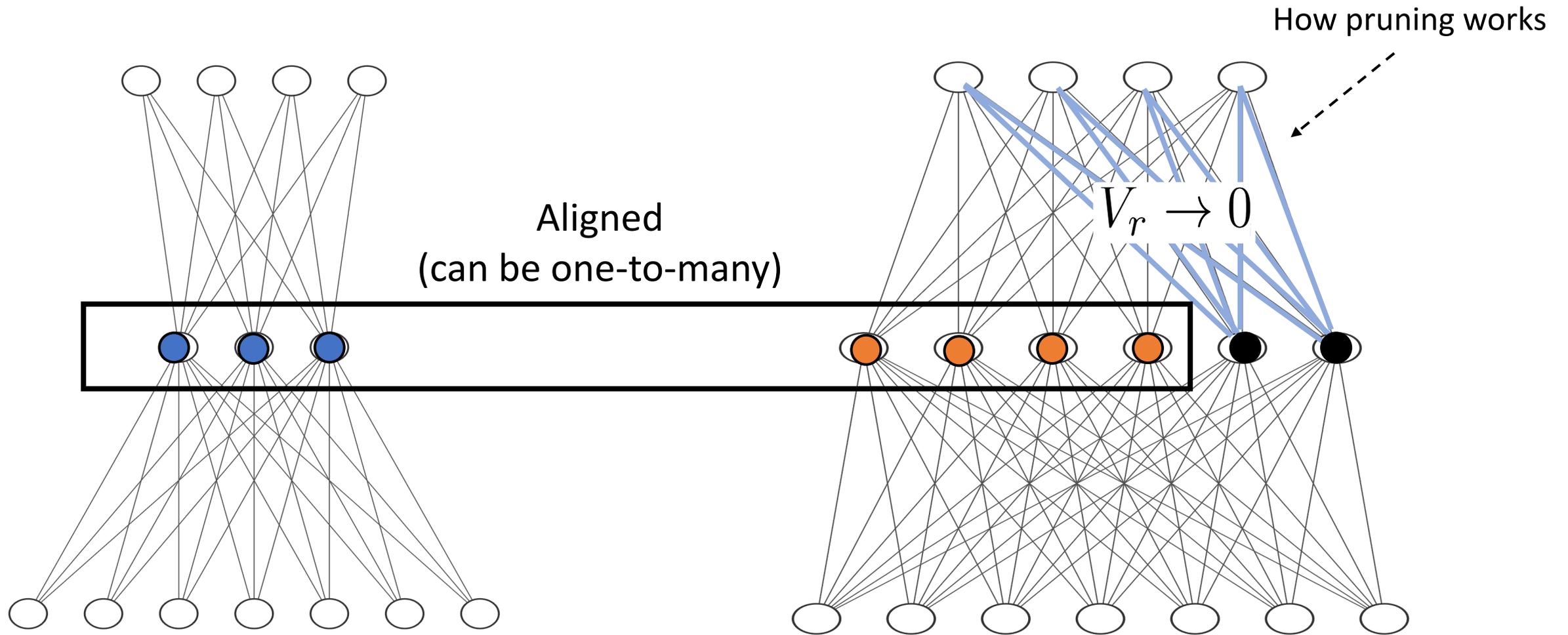
$\mathbf{g}_1(x) = \mathbf{0}$  for all  $x \in R_0$   
(all input gradients at layer 1 is  
*zero* at all training samples)

Teacher node  $j$  is **observed**  
by a student node  $k$

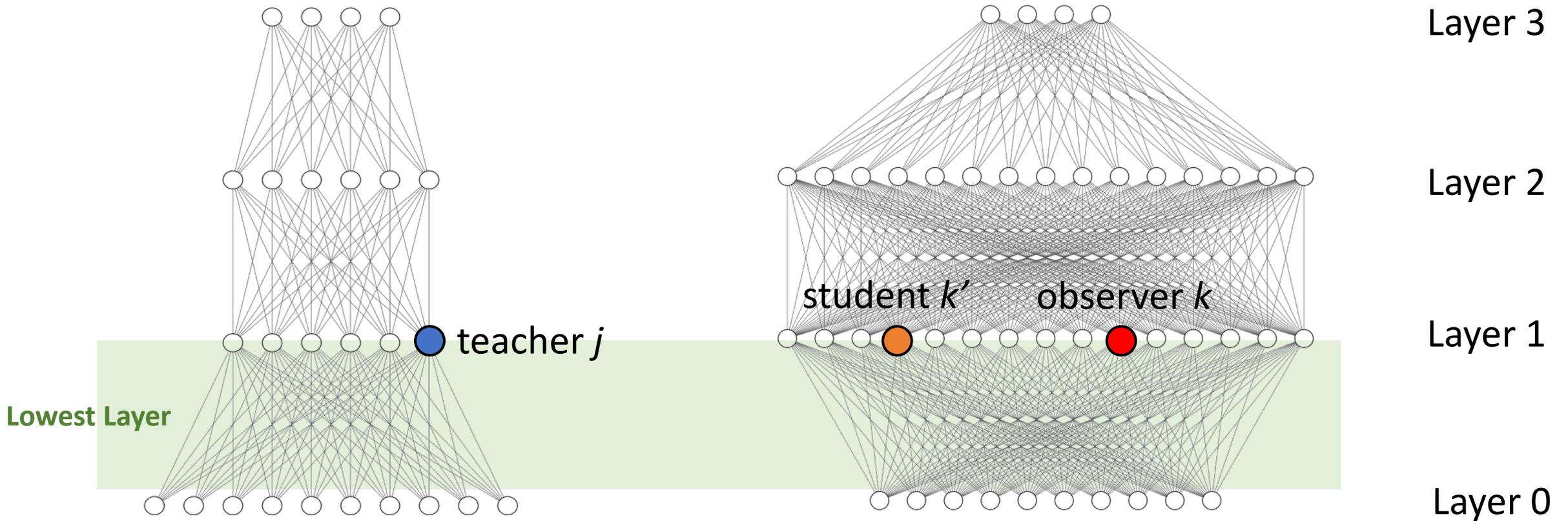


Teacher  $j$  is **aligned with**  
at least one student  $k'$

# What happens to unaligned students?



# Multi-Layer case: Alignment could happen!



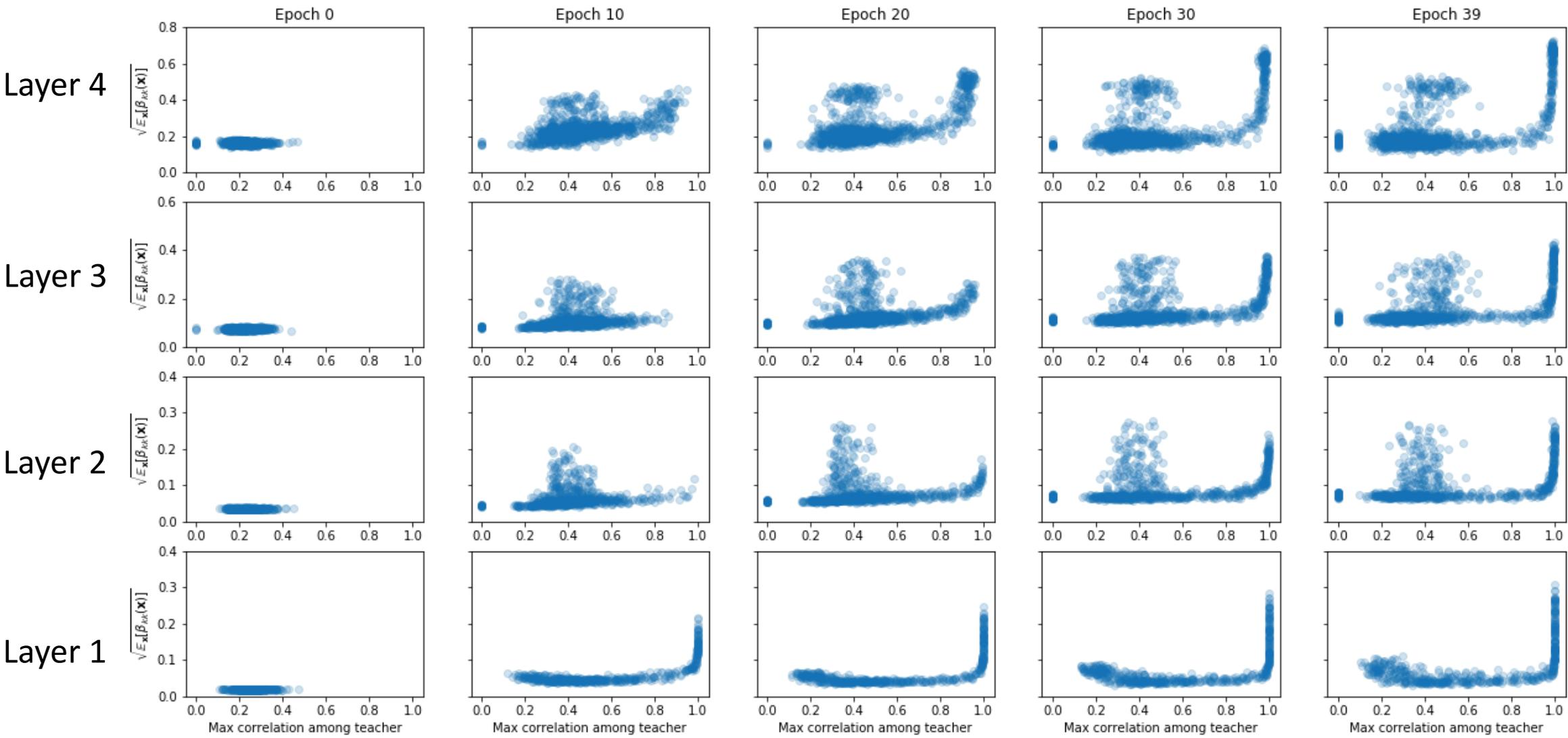
$$\alpha_k^T(\mathbf{x})\mathbf{f}^*(\mathbf{x}) - \beta_k^T(\mathbf{x})\mathbf{f}(\mathbf{x}) = \mathbf{0}$$

Piece wise constant, apply the same logic **per region!**

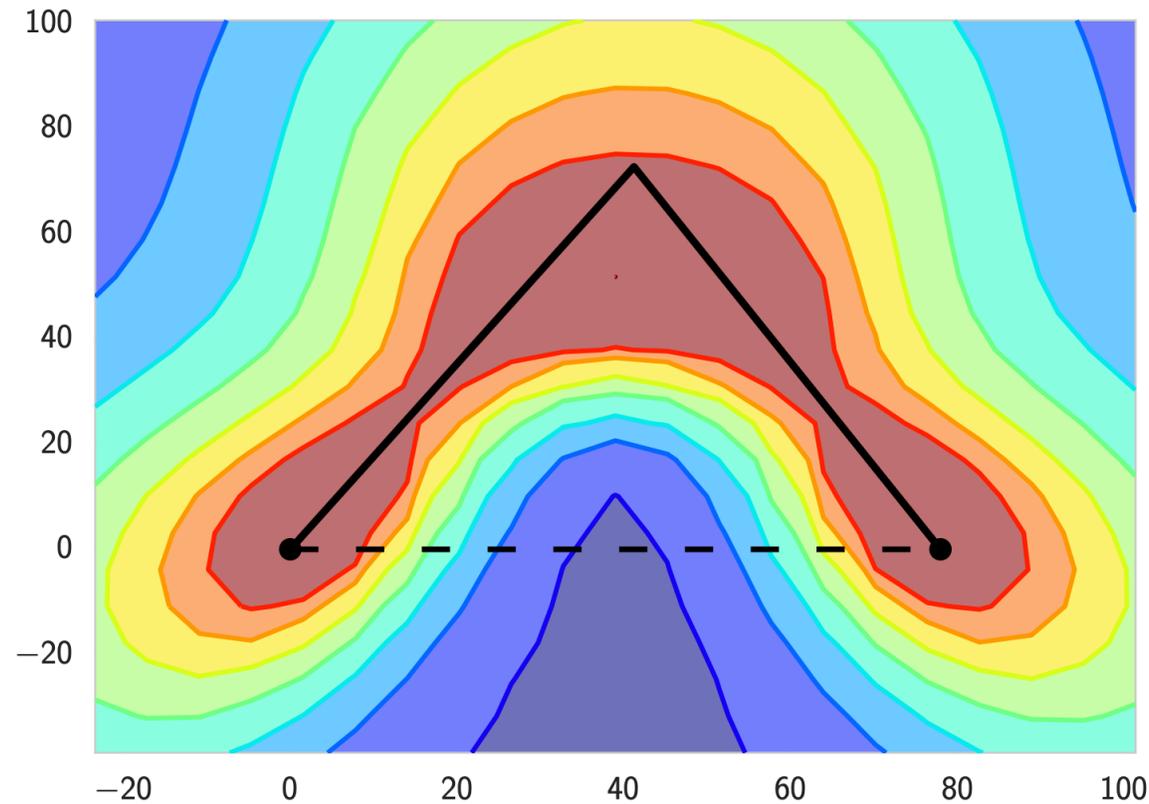
For 2-layer:

$$\sqrt{\mathbb{E}_{\mathbf{x}} [\beta_{kk}(\mathbf{x})]} = \|\mathbf{v}_k\|$$

# Training Progresses



# Solutions can be connected by line segments

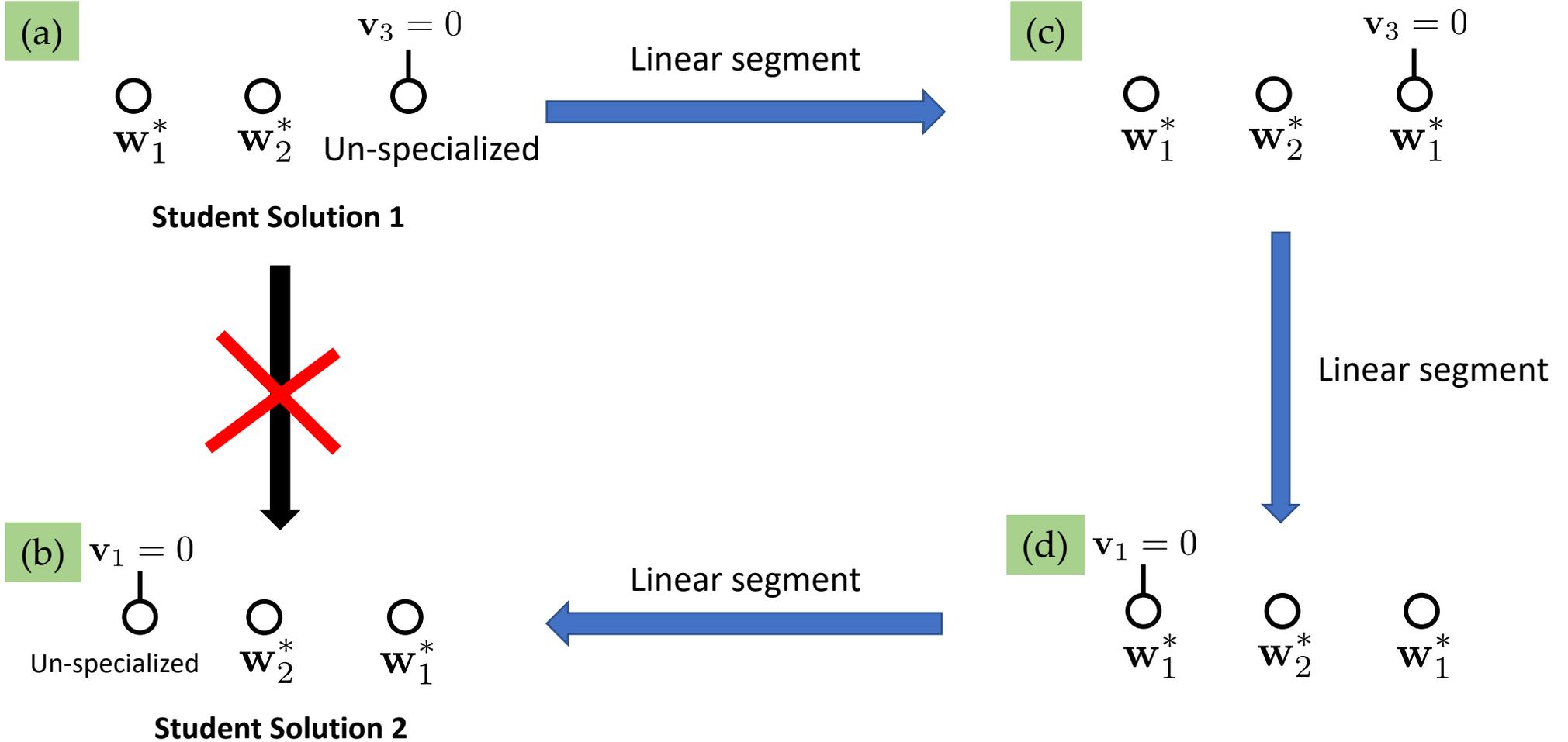


*[Loss Surfaces, Mode Connectivity, and Fast Ensembling of DNNs, Garipov et al. NeurIPS 2018]*

*[Essentially No Barriers in Neural Network Energy Landscape, Draxler et al, 2018]*

*[Explaining Landscape Connectivity of Low-cost Solutions for Multilayer Nets, Kuditipudi et al, 2019]*

# Our Explanation



# Realistic Case

Student Specialization with **2-layers** ReLU nets,  
**Small** Gradient and **Finite** Samples

# Polynomial Complexity for 2-layered Network

To achieve  $\epsilon$ -alignment between a teacher  $j$  and student  $k$

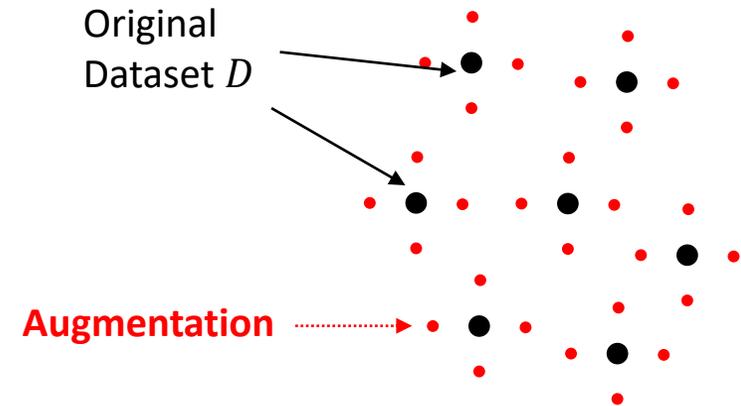
$$K_1 = m_1 + n_1$$

Small gradient

$$\|\mathbf{g}_1(\mathbf{x}, \hat{\mathcal{W}})\|_\infty \leq \frac{\alpha_{kj}}{5K_1^{3/2}\sqrt{d}}\epsilon, \mathbf{x} \in D'$$

Sample Complexity of original Dataset  $D$

$$N = \mathcal{O}(K_1^{5/2}d^2\epsilon^{-1}\kappa^{-1})$$

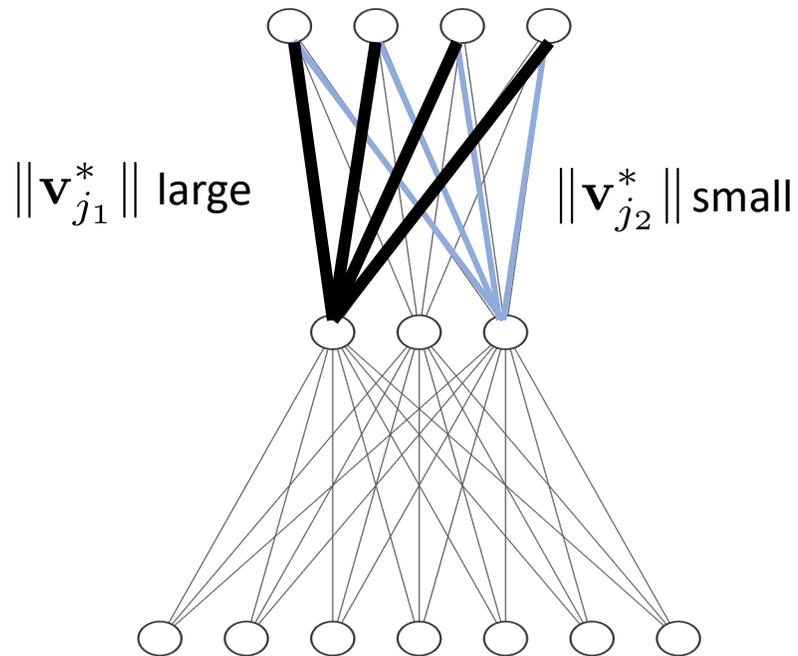


**Teacher-agnostic** augmentation

$$D' = \text{Aug}(D)$$

$$|D'| = (2d+1)|D|$$

# Lesson 1: Stronger teacher node learns faster



Gradient Condition:

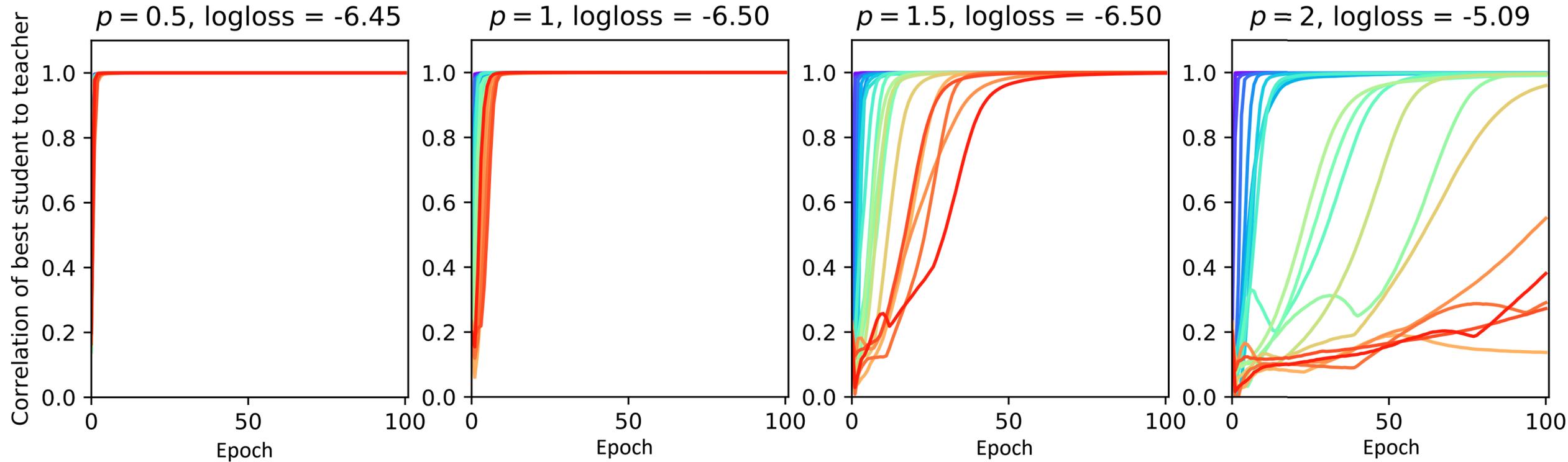
$$\alpha_{kj} := \mathbf{v}_k^T \mathbf{v}_j^*$$
$$\|\mathbf{g}_1(\mathbf{x}, \hat{\mathcal{W}})\|_\infty \leq \frac{\alpha_{kj}}{5K_1^{3/2}\sqrt{d}} \epsilon, \mathbf{x} \in D'$$

Strong teacher nodes are learned faster

1. Robust to Noise! 😊
2. Hard to learn weak teacher nodes 😞

# Weak teacher nodes are slow to learn

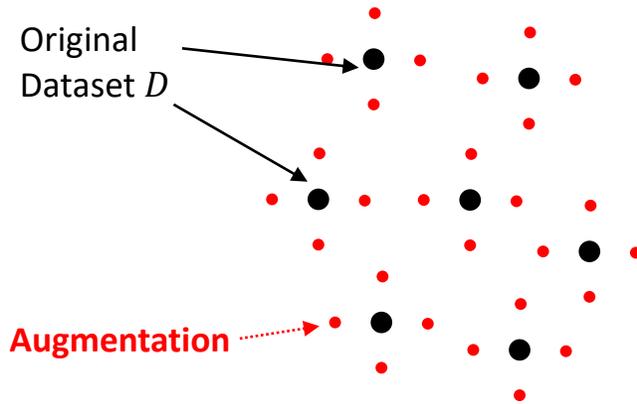
Teacher  $j$ :  
 $\|\mathbf{v}_j^*\| \propto 1/j^p$



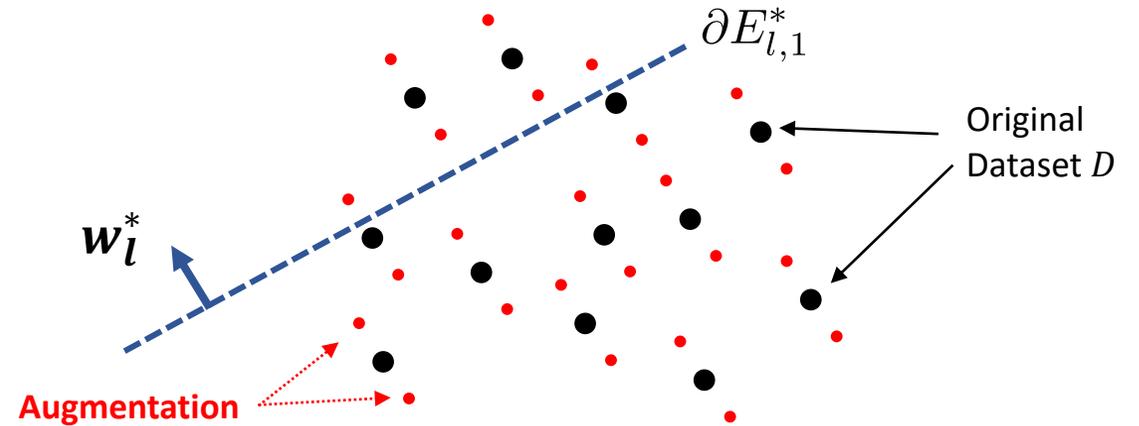
Weak teacher  Strong teacher

# Lesson 2: Data augmentation Matters

## Teacher-agnostic augmentation



## Teacher-aware augmentation



Small gradient

$$\|g_1(\mathbf{x}, \hat{W})\|_\infty \leq \frac{\alpha_{kj}}{5K_1^{3/2}\sqrt{d}}\epsilon$$

$$\|g_1(\mathbf{x}, \hat{W})\|_\infty \leq \frac{\alpha_{kj}}{5K_1^{3/2}}\epsilon$$

**No  $\sqrt{d}$**

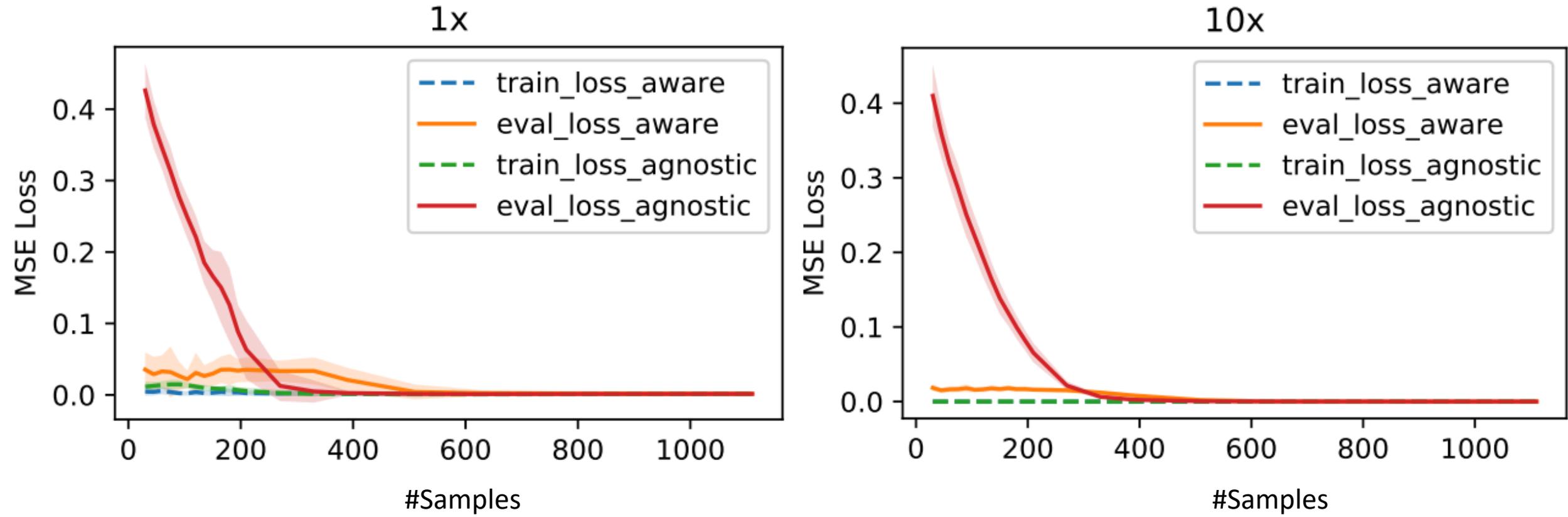
Sample Complexity of original Dataset  $D$

$$N = \mathcal{O}(K_1^{5/2}d^2\epsilon^{-1}\kappa^{-1})$$

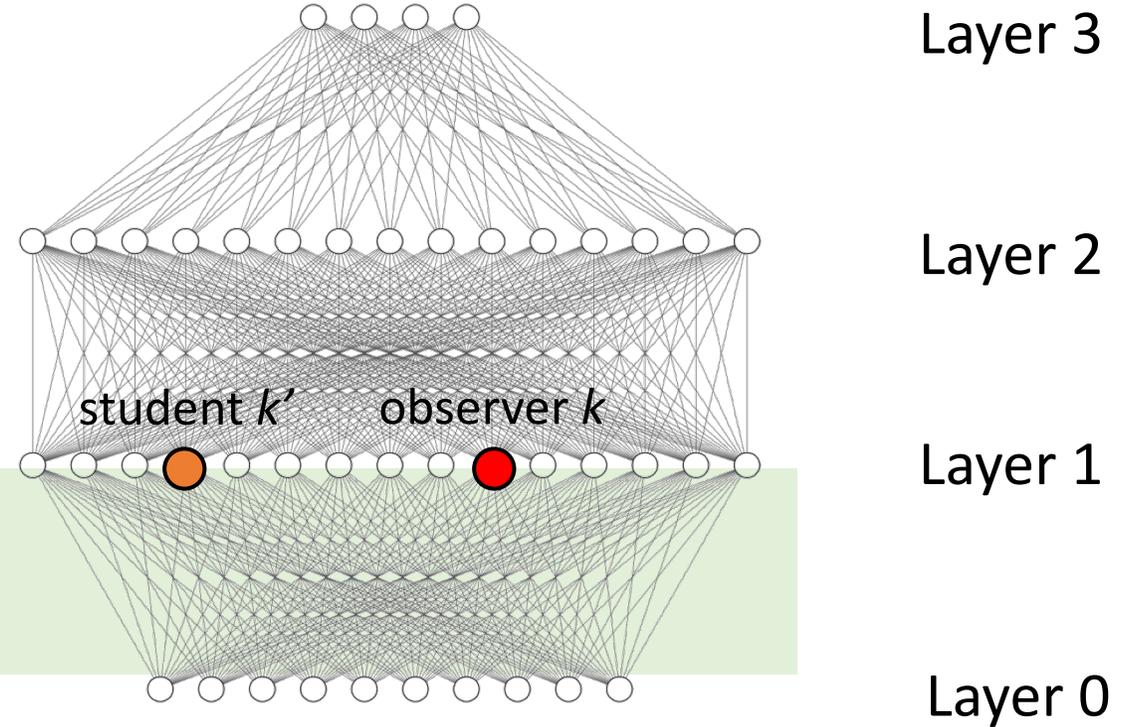
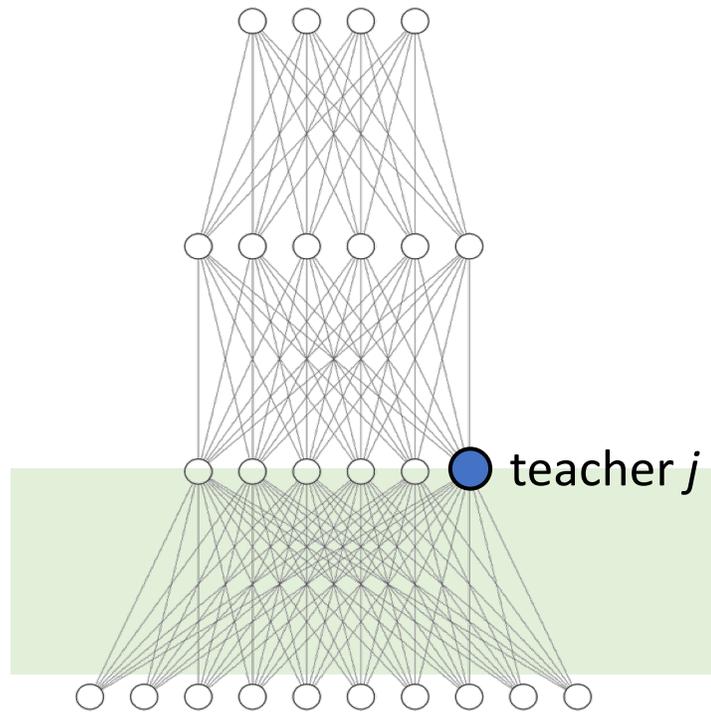
$$\mathcal{O}(K_1^{5/2}d\epsilon^{-1}\kappa^{-1})$$

**Linear w.r.t  $d$**

# Teacher-Agnostic versus Teacher-aware



# Multi-layer case



Small gradient

$$\|\mathbf{g}_1(\mathbf{x}, \hat{\mathcal{W}})\|_\infty \leq \frac{\min_{R \in \mathcal{R}} \alpha_{kj}(R)}{5Q^{3/2}\sqrt{d}} \epsilon, \text{ for } \mathbf{x} \in D'$$

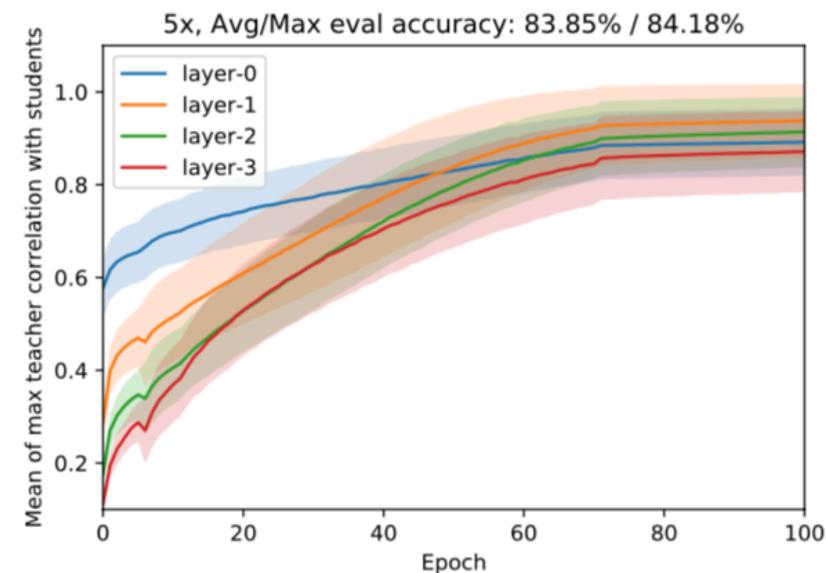
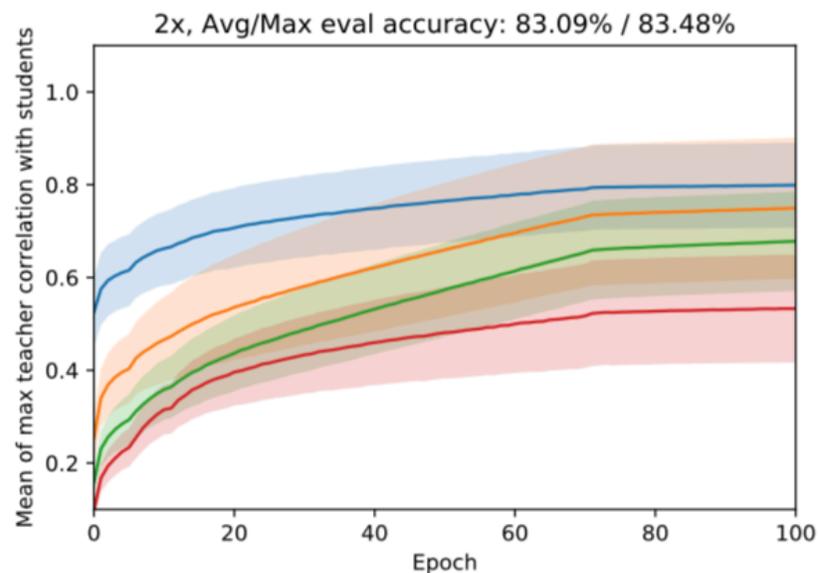
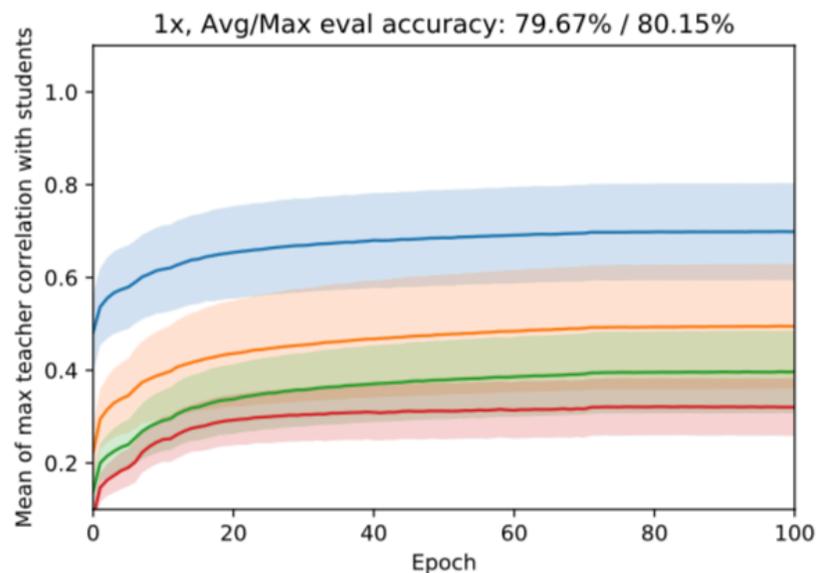
Sample Complexity of original Dataset  $D$

$$\mathcal{O}(Q^{5/2}d^2\epsilon^{-1}\kappa^{-1})$$

# CIFAR 10

1. Train a conv teacher network of size 64-64-64-64.
2. **[Construct Oracle]** Prune the teacher network to [45-32-32-20]
3. Then train a student network to mimic teacher's output (before softmax)

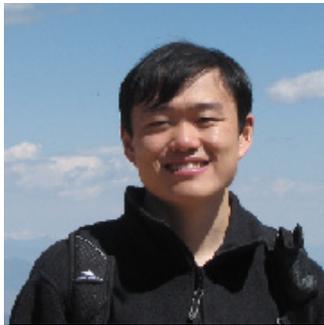
The student network has more parameters



A theoretical framework that explains

1. Why **self-supervised learning** with deep ReLU models works
2. Why a good representation is learned without supervision
3. Why BYOL doesn't need negative samples

# Understanding Self-supervised Learning with Dual Deep Networks



Yuandong Tian



Lantao Yu

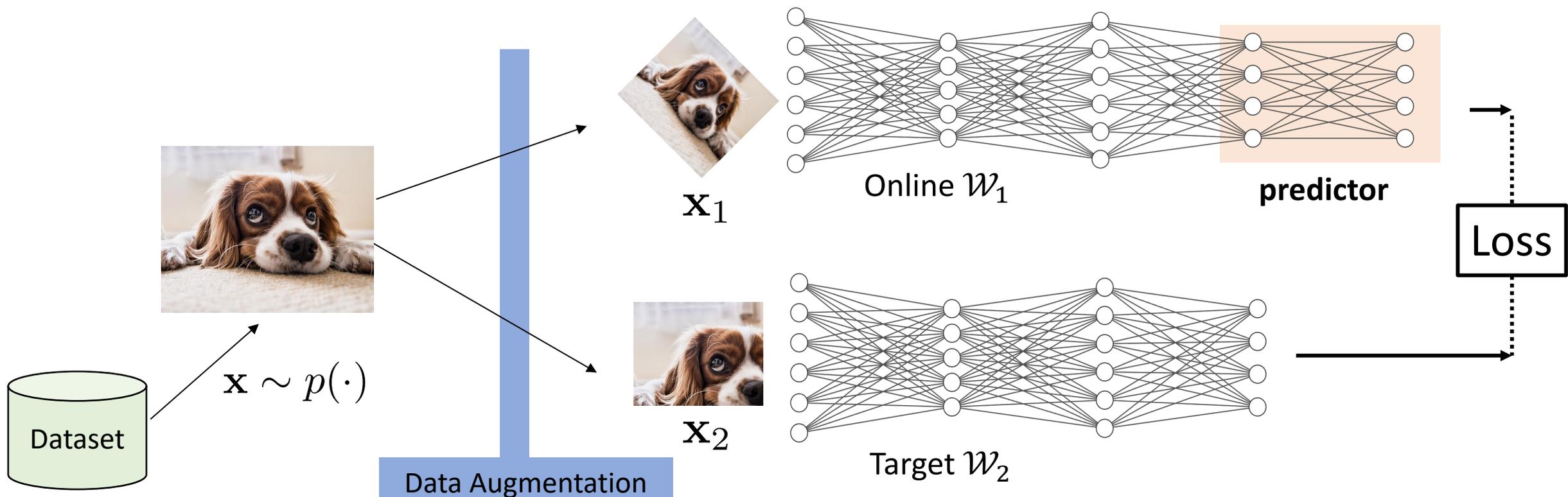


Xinlei Chen



Surya Ganguli

# Self-supervised Learning (SSL)



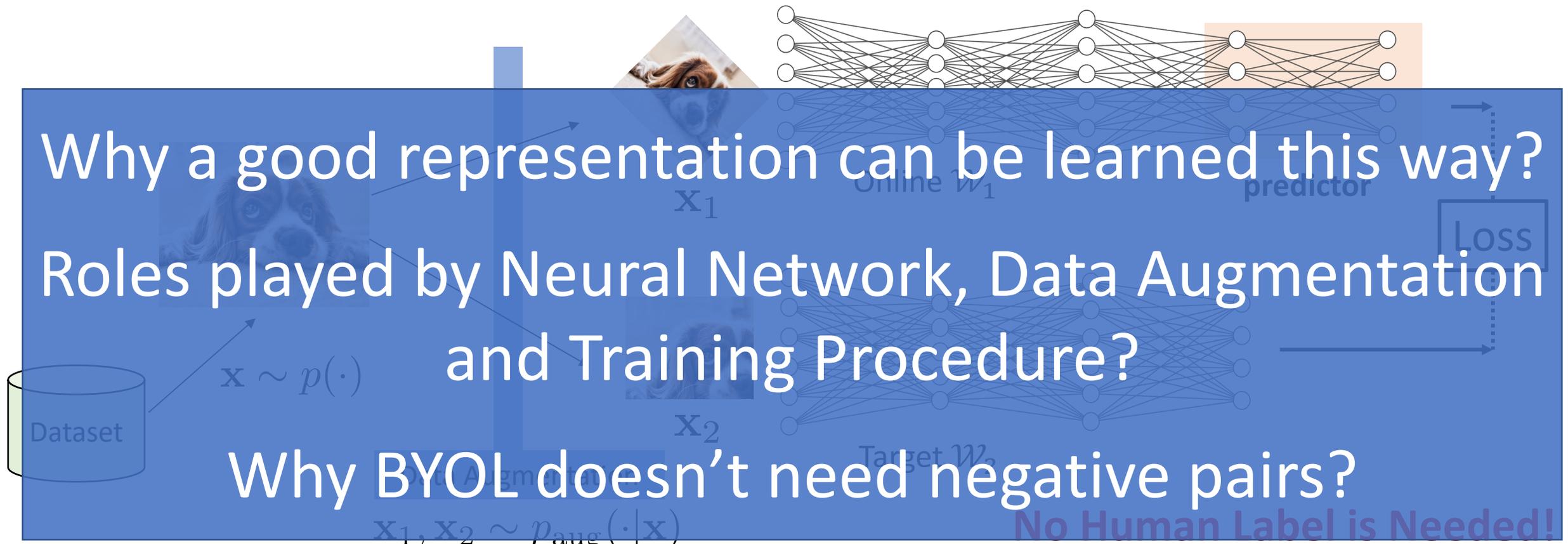
$$\mathbf{x}_1, \mathbf{x}_2 \sim p_{\text{aug}}(\cdot | \mathbf{x})$$

**No Human Label is Needed!**

**BYOL:** [J. Grill, *Bootstrap your own latent: A new approach to self-supervised Learning*, arXiv 2020]

**SimCLR:** [T. Chen, *A Simple Framework for Contrastive Learning of Visual Representations*, ICML 2020]

# Self-supervised Learning (SSL)



Why a good representation can be learned this way?  
Roles played by Neural Network, Data Augmentation  
and Training Procedure?

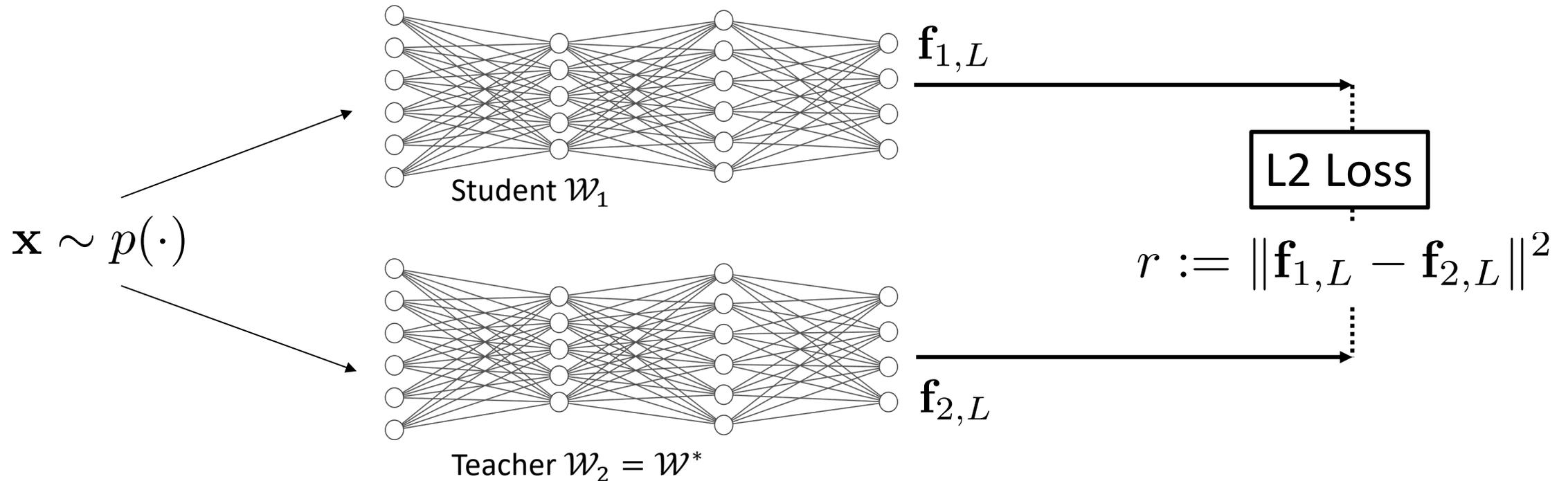
Why BYOL doesn't need negative pairs?

No Human Label is Needed!

**BYOL:** [J. Grill, Bootstrap your own latent: A new approach to self-supervised Learning, arXiv 2020]

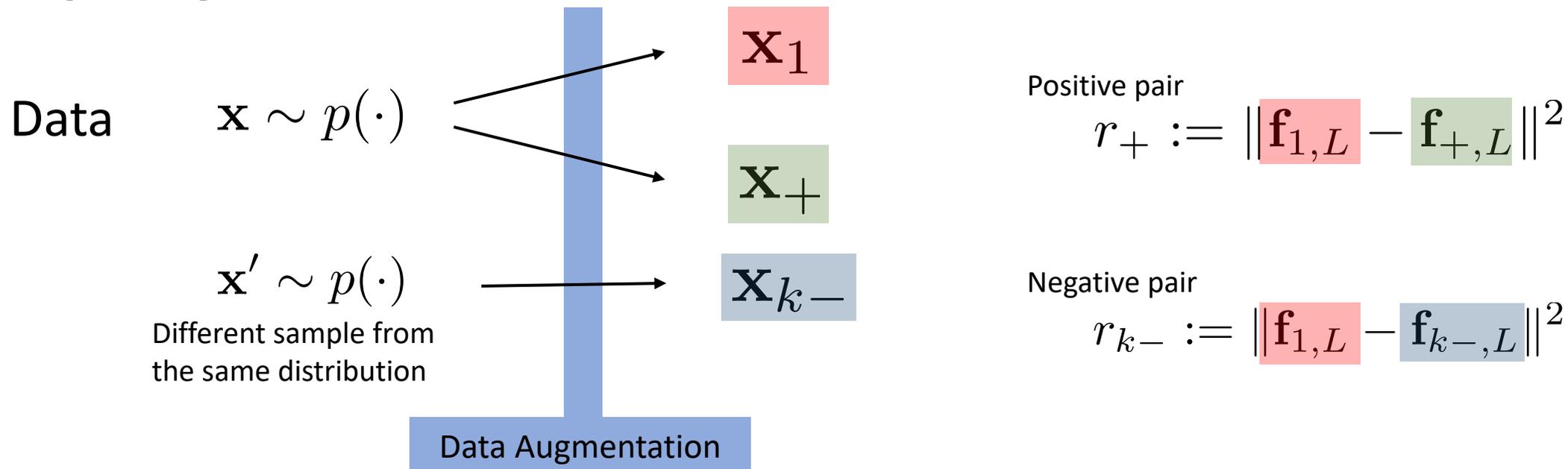
**SimCLR:** [T. Chen, A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020]

# Similarity with Teacher Student Setting



**The mathematical framework is similar!**

# SimCLR

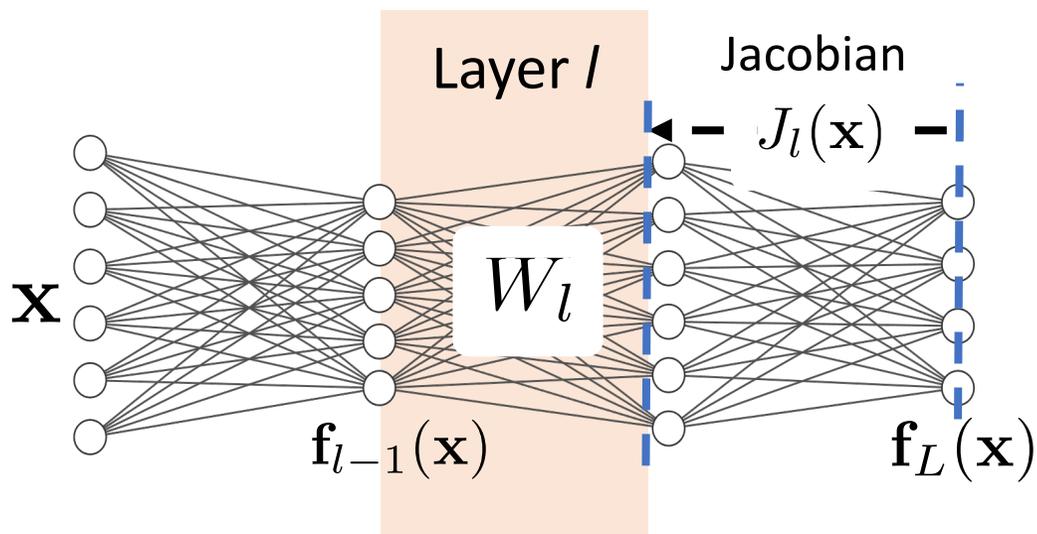


## InfoNCE

$$L(r_+, r_{1-}, r_{2-}, \dots, r_{K-}) := -\log \frac{e^{-r_+/\tau}}{e^{-r_+/\tau} + \sum_{k=1}^H e^{-r_{k-}/\tau}}$$

If  $|u| = |v| = 1$ , then the formulation is the same as SimCLR's formulation  
Since  $-r = -|u - v|^2 = 2\text{sim}(u, v) - 2$

# The Covariance Operator



Weight Update for SimCLR at layer  $l$ :

$$W_l(t + 1) = W_l(t) + \alpha \Delta W_l(t)$$

Learning rate

Connection

$$K_l(\mathbf{x}) := \mathbf{f}_{l-1}(\mathbf{x}) \otimes J_l^\top(\mathbf{x})$$

Augment-mean Connection

$$\bar{K}_l(\mathbf{x}) := \mathbb{E}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot | \mathbf{x})} [K_l(\mathbf{x}')] ]$$

**Covariance operator (PSD)**

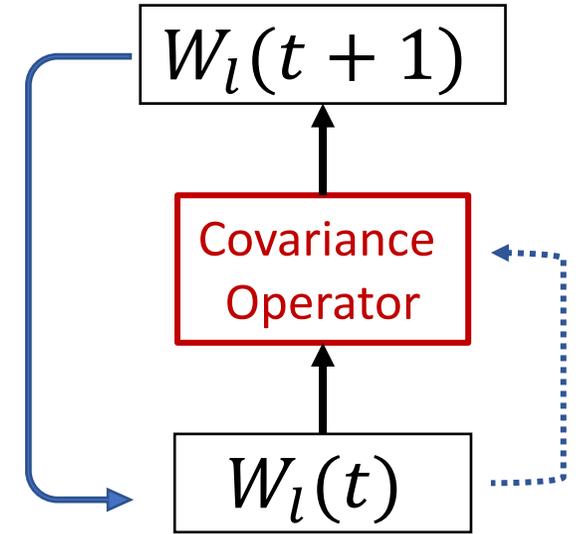
$$\text{vec}(\Delta W_l(t)) = \beta \mathbb{V}_x [\bar{K}_l(\mathbf{x})] \text{vec}(W_l(t))$$

Positive number related to Contrastive loss

# What does it mean?

**The Covariance Operator**  $\mathbb{V}_{\mathbf{x}}[\bar{K}_l(\mathbf{x}; \mathcal{W}(t))]$

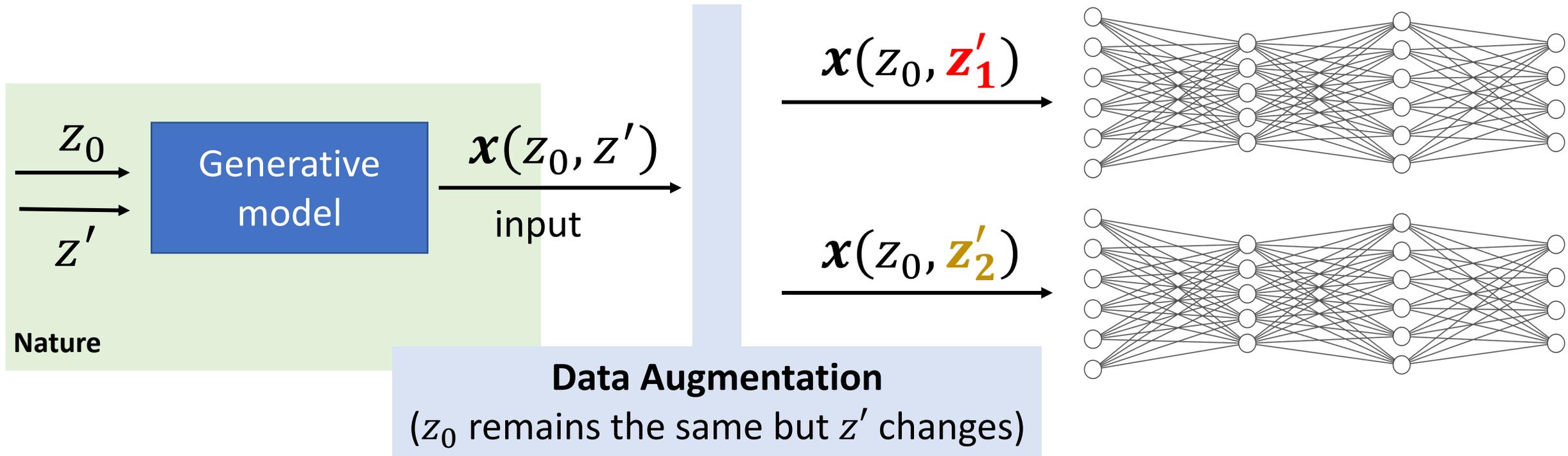
- Always PSD at any stage of training
- Weight at each layer undergoes a PSD transformation
- Strong eigen mode leads strong weight growth along that direction



## What are the strong eigen models in the covariance operator?

**To understand that, we need a generative model of the data.**

# Using Generative Models to understand Covariance Operator

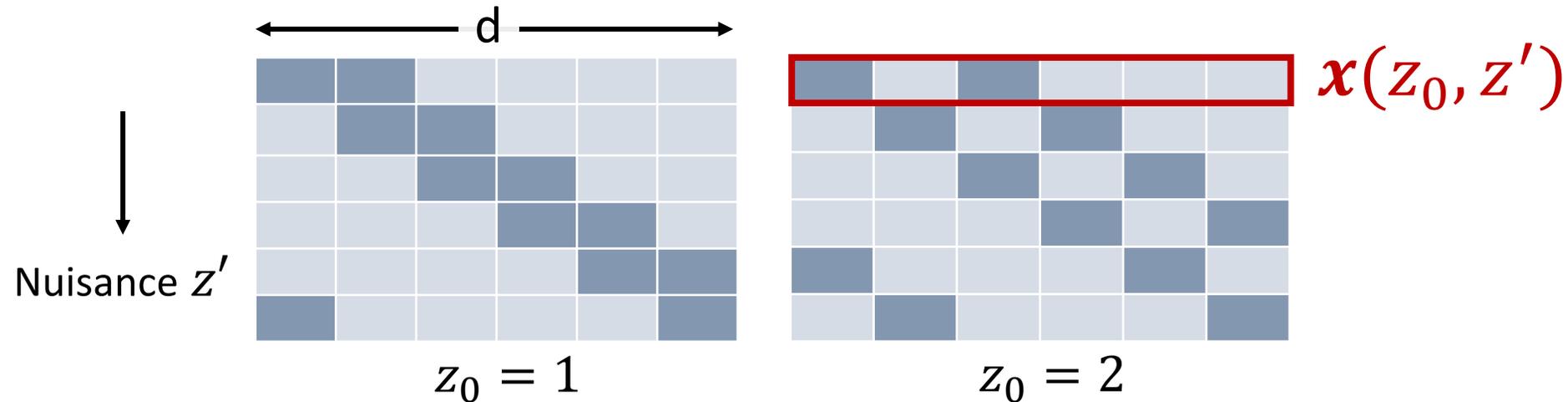


$z_0$ : Class (sample) label

$z'$ : Nuisance Transformations given by Data Augmentation

# One-layer one-neuron example

Two objects **11** and **101** translating in 1D space



$$\mathbb{V}_{z_0} [\bar{K}(z_0)] = \frac{1}{4d^2} \mathbf{u}\mathbf{u}^\top$$

$$\mathbf{u} := \mathbf{x}_{11} + \mathbf{x}_{00} - \mathbf{x}_{01} - \mathbf{x}_{10}$$

Feature to represent pattern 10

Linear neuron: Nothing is learned.

ReLU neuron: Enforce what is initialized!

# A two-layer example

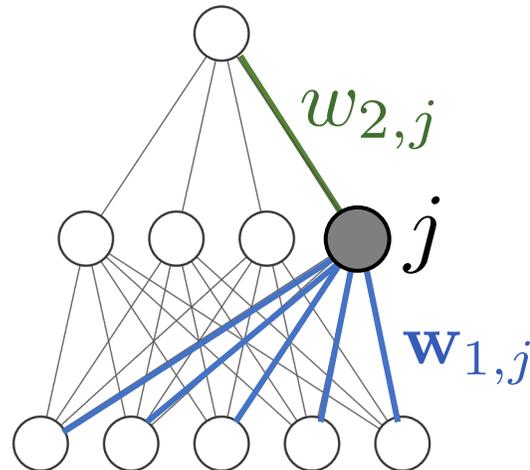
Augment-Average Connection for both layers:

$$\bar{K}_1(z) = [w_{2,1} \mathbf{u}_1, \dots, w_{2,n_1} \mathbf{u}_{n_1}] \quad \bar{K}_2(z) = [\mathbf{w}_{1,1}^\top \mathbf{u}_1, \dots, \mathbf{w}_{1,n_1}^\top \mathbf{u}_{n_1}]$$

$$\text{where } \mathbf{u}_j(z) := \mathbb{E}_{z'|z} [\mathbf{x}(z, z') \mathbb{I}(\mathbf{w}_{1,j}^\top \mathbf{x}(z, z') \geq 0)]$$

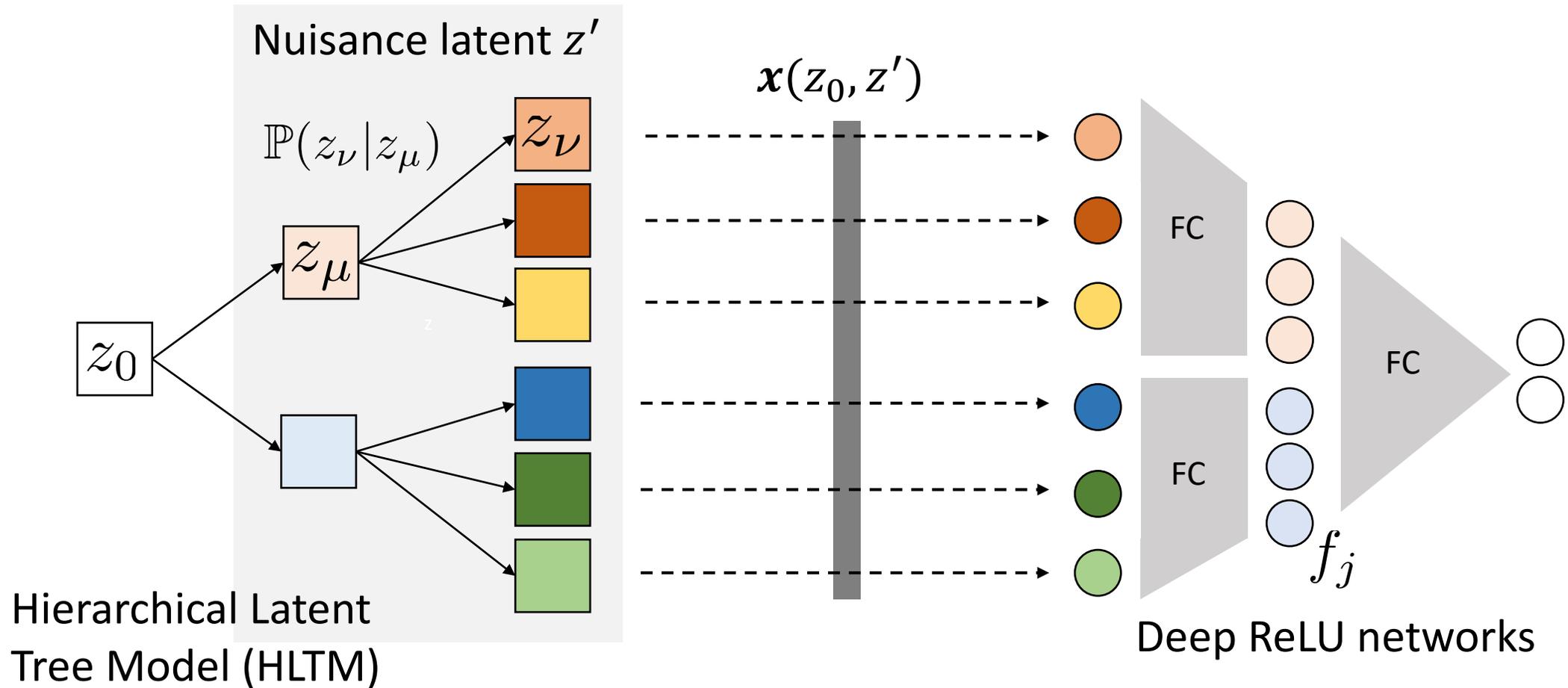
**Theorem 4.** *If  $\text{Cov}_z [\mathbf{u}_j, \mathbf{u}_k] = 0$  for  $j \neq k$ , then the time derivative of  $w_{2,j}$  and  $\mathbf{w}_{1,j}$  satisfies:*

$$\dot{w}_{2,j} = w_{2,j} \mathbf{w}_{1,j}^\top A_j \mathbf{w}_{1,j}, \quad \dot{\mathbf{w}}_{1,j} = w_{2,j}^2 A_j \mathbf{w}_{1,j}, \quad \text{where } A_j := \mathbb{V}_z[\mathbf{u}_j(z)].$$

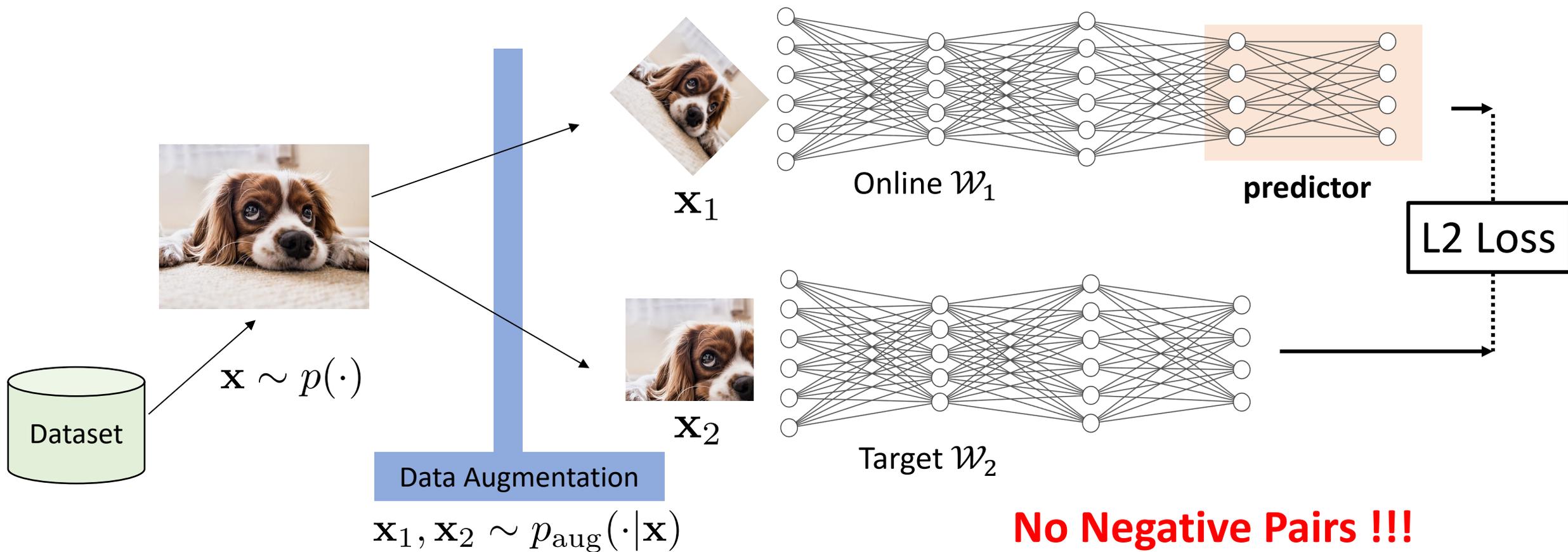


Weights of two layer are **enforcing** each other

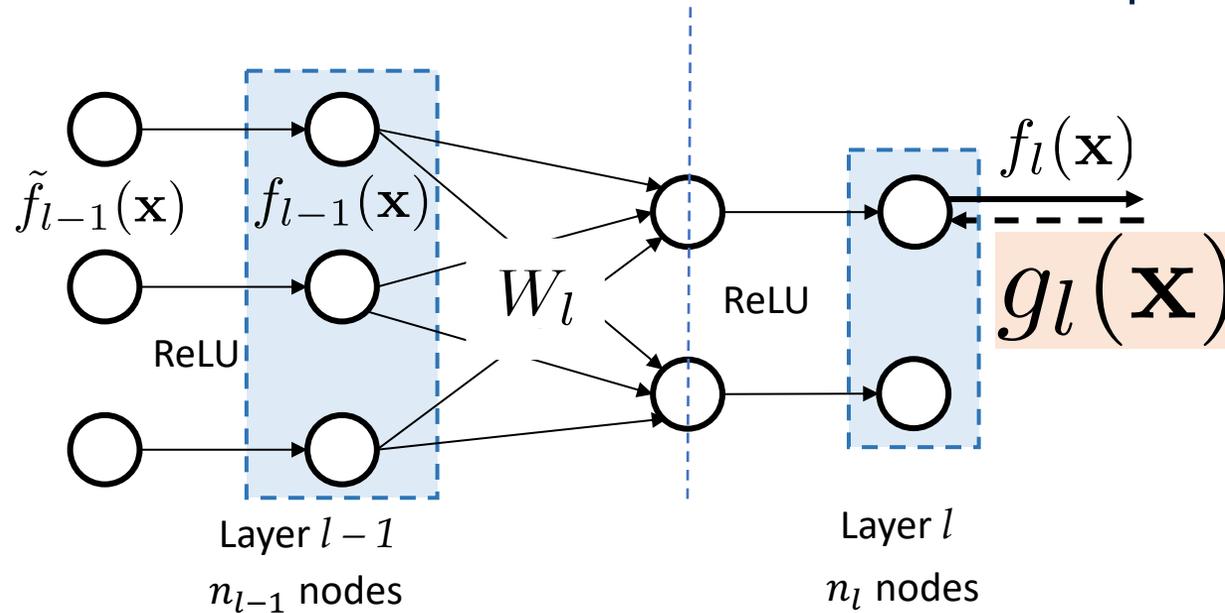
# Hierarchical Latent Tree Models (HLTM)



# BYOL Setting



# Is BatchNorm the Secret Weapon?



## Zero-mean property.

After BN, Backpropagated Gradient is zero-mean in each minibatch:

$$\tilde{\mathbf{g}}_l^i := \mathbf{g}_l^i - \frac{1}{|B|} \sum_{i \in B} \mathbf{g}_l^i = \mathbf{g}_l^i - \bar{\mathbf{g}}_l$$

With **Normalization** and  $\mathcal{W} \neq \mathcal{W}'$

$$\begin{aligned}
 \text{vec} \left( \widetilde{\Delta W_l} \right) &= \text{vec}(\Delta W_l) + \text{vec}(\delta W_l^{\text{BN}}) \\
 &= \text{vec}(\Delta W_l)_{\text{sym}} \quad \text{Enforce similar representation within data augmentation} \\
 &\quad - \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) + \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W'_l)
 \end{aligned}$$

$$\text{vec}(\Delta W_l)_{\text{sym}} = -\mathbb{E}_{\mathbf{x}} \{ \mathbb{V}_{\mathbf{x}'} [K_l(\mathbf{x}')] \} \text{vec}(W_l)$$

# Why can BYOL work?

If  $\mathcal{W}$  has an extra predictor  $\left\| -\nabla_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \right\| \ll \left\| \text{Cov}_{\mathbf{x}}[\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \right\|$

$$\begin{aligned} \text{vec} \left( \widetilde{\Delta W}_l \right) &= \text{vec}(\Delta W_l) + \text{vec}(\delta W_l^{\text{BN}}) \\ &= \text{vec}(\Delta W_l)_{\text{sym}} \\ &\quad - \nabla_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) + \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \end{aligned}$$

**Negated** covariance operator  
(but second order w.r.t the predictor!!)

**Approximate covariance operator**

# BYOL Setting (Top-1 Performance in STL-10)

**No predictor, things do not work**

-	EMA	BN	EMA, BN
$38.7 \pm 0.6$	$39.3 \pm 0.9$	$33.0 \pm 0.3$	$32.8 \pm 0.5$

**BN is critical**

**EMA is helping**

P	P, EMA	P, BN	P, EMA, BN
$39.5 \pm 3.1$	$44.4 \pm 3.2$	$63.6 \pm 1.06$	<b><math>78.1 \pm 0.3</math></b>

# Zero-mean Gradient matters.

## Ablation Study of Batch components

-	$\mu$	$\sigma$	$\mu, \sigma$	$\mu^\#$
$43.9 \pm 4.2$	$64.8 \pm 0.6$	$72.2 \pm 0.9$	<b><math>78.1 \pm 0.3</math></b>	$44.2 \pm 7.0$
$\sigma^\#$	$\mu^\#, \sigma$	$\mu, \sigma^\#$	$\mu^\#, \sigma^\#$	
$54.2 \pm 0.6$	<b><math>48.3 \pm 2.7</math></b>	$76.3 \pm 0.4$	<b><math>47.0 \pm 8.1</math></b>	

$\mu$   $x = x - x.\text{mean}(0)$

$\sigma$   $x = x / x.\text{std}(0)$

$\mu^\#$   $x = x - x.\text{mean}(0).\text{detach}()$

$\sigma^\#$   $x = x / x.\text{std}(0).\text{detach}()$

# Reinitializing Predictors Works

Table 5: Top-1 performance of BYOL using reinitialization of the predictor every  $T$  epochs.

	Original BYOL	ReInit $T = 5$	ReInit $T = 10$	ReInit $T = 20$
STL-10 (100 epochs)	78.1	78.6	<b>79.1</b>	79.0
ImageNet (60 epochs)	60.9	61.9	<b>62.4</b>	<b>62.4</b>

The predictor is not necessarily “optimal” as suggested in the original BYOL paper.

# Understanding Adversarial Samples using Teacher-Student



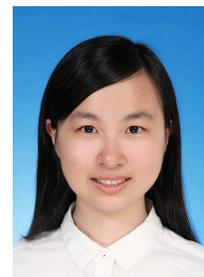
Zhoulin Yang\*



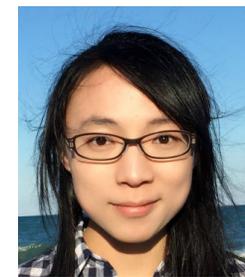
Zhaoxi Chen



Tiffany Cai



Xinyun Chen



Bo Li



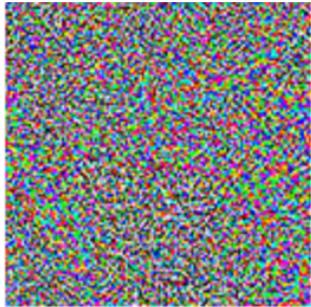
Yuandong Tian \*

\* = Equal Contribution

# Adversarial Samples



+  $\epsilon$



=

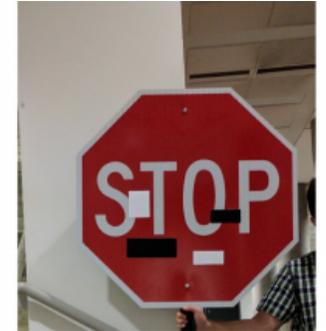


"panda"

57.7% confidence

"gibbon"

99.3% confidence



Stop sign  $\rightarrow$  a 45 mph sign



Adversarial Anti-Facial Recognition...

**\$20.23**

Redbubble

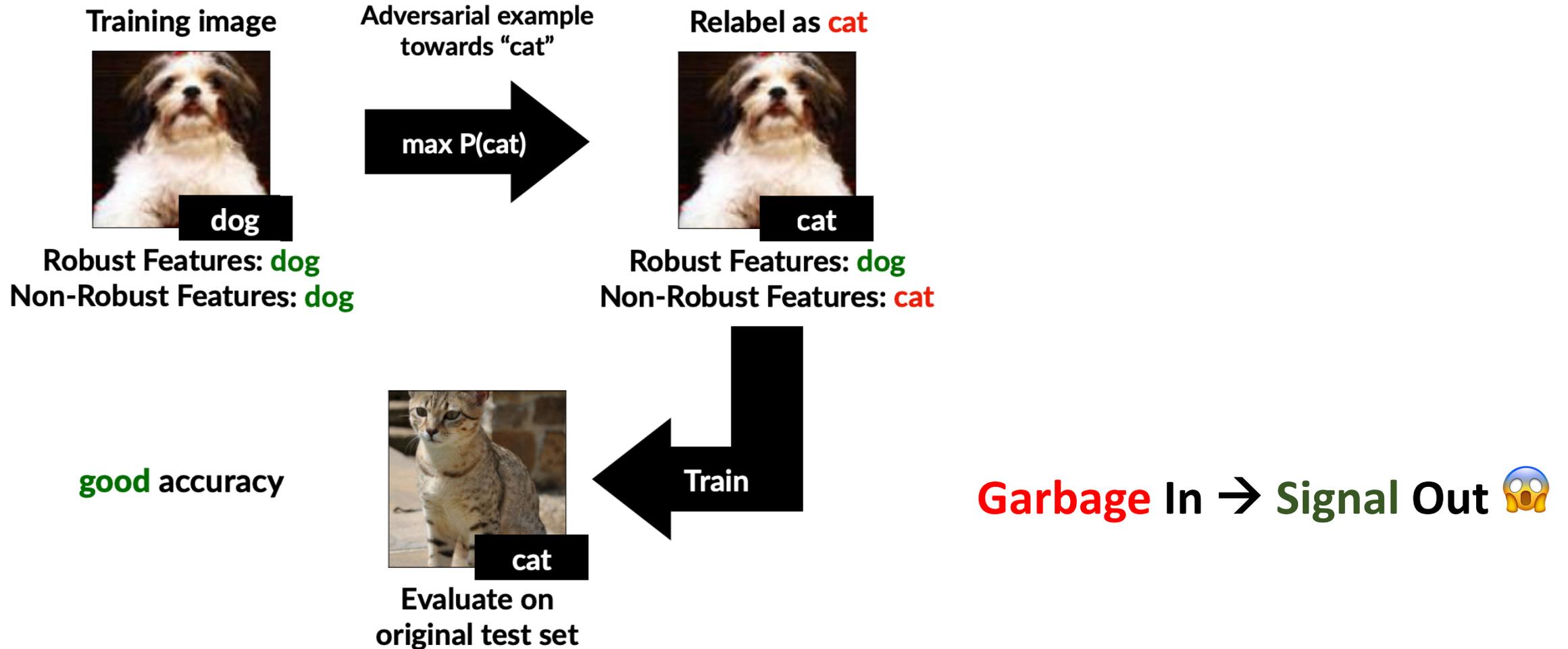


Adversarial Anti-Facial Recognition...

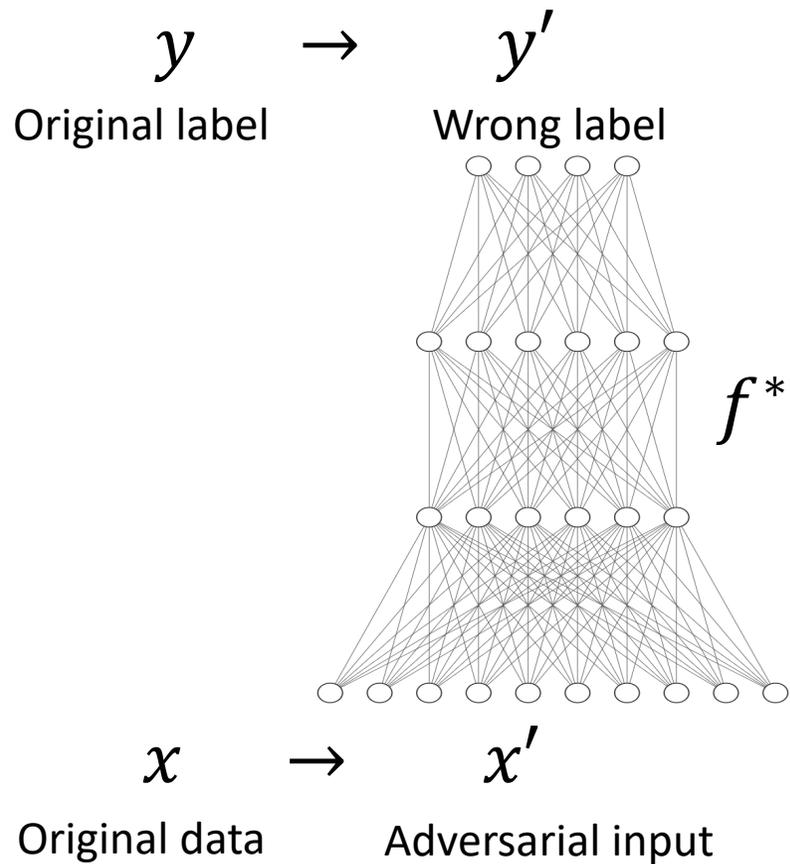
**\$28.69**

Redbubble

# Adversarial Samples: not Bugs but Features

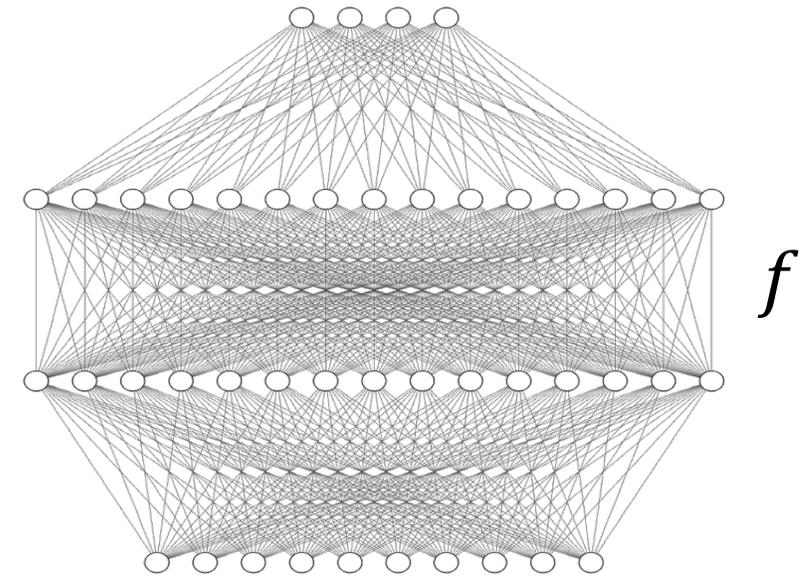


# A Natural Explanation



**Normal training:** train with  $(x, y)$

**Fancy training:** train with  $(x', y' = f^*(x'))$   
(wrong sample, correct mapping)

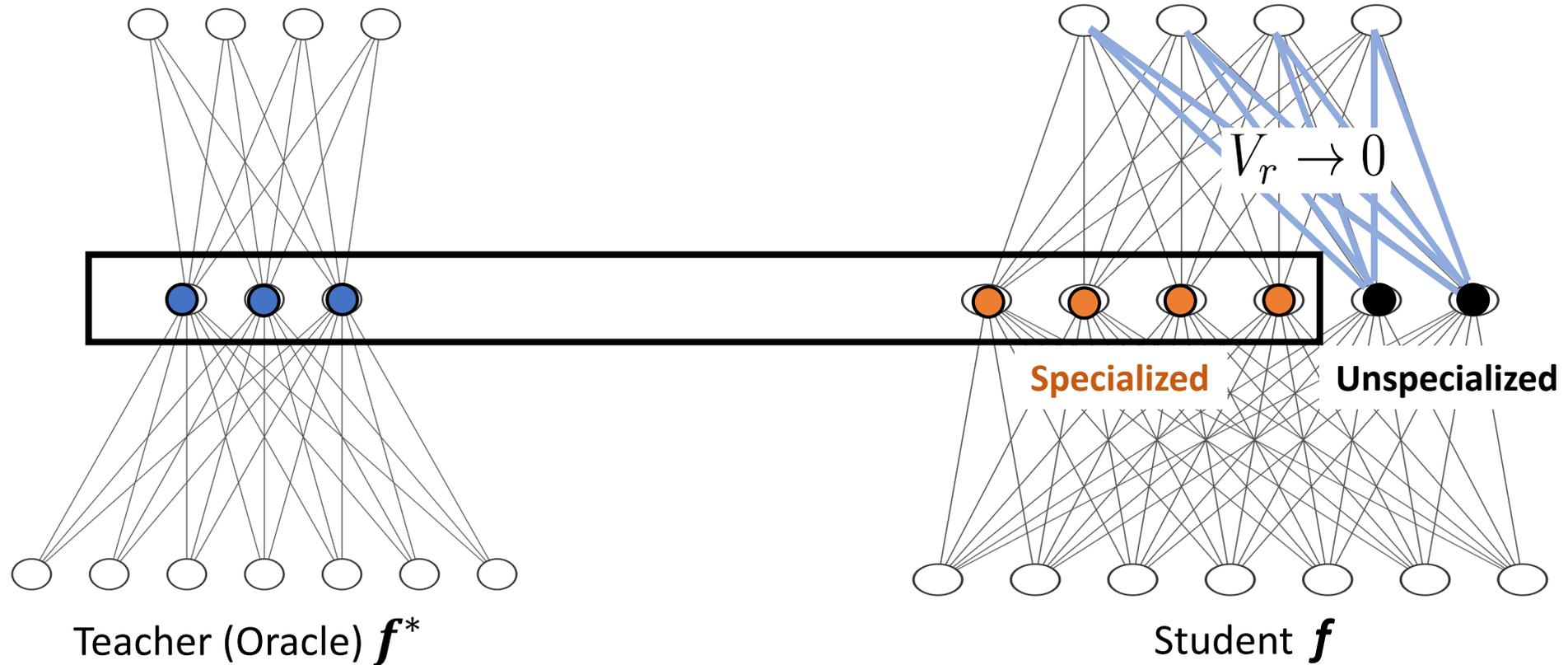


Both are valid pairs from the original network  
→ They both train well.

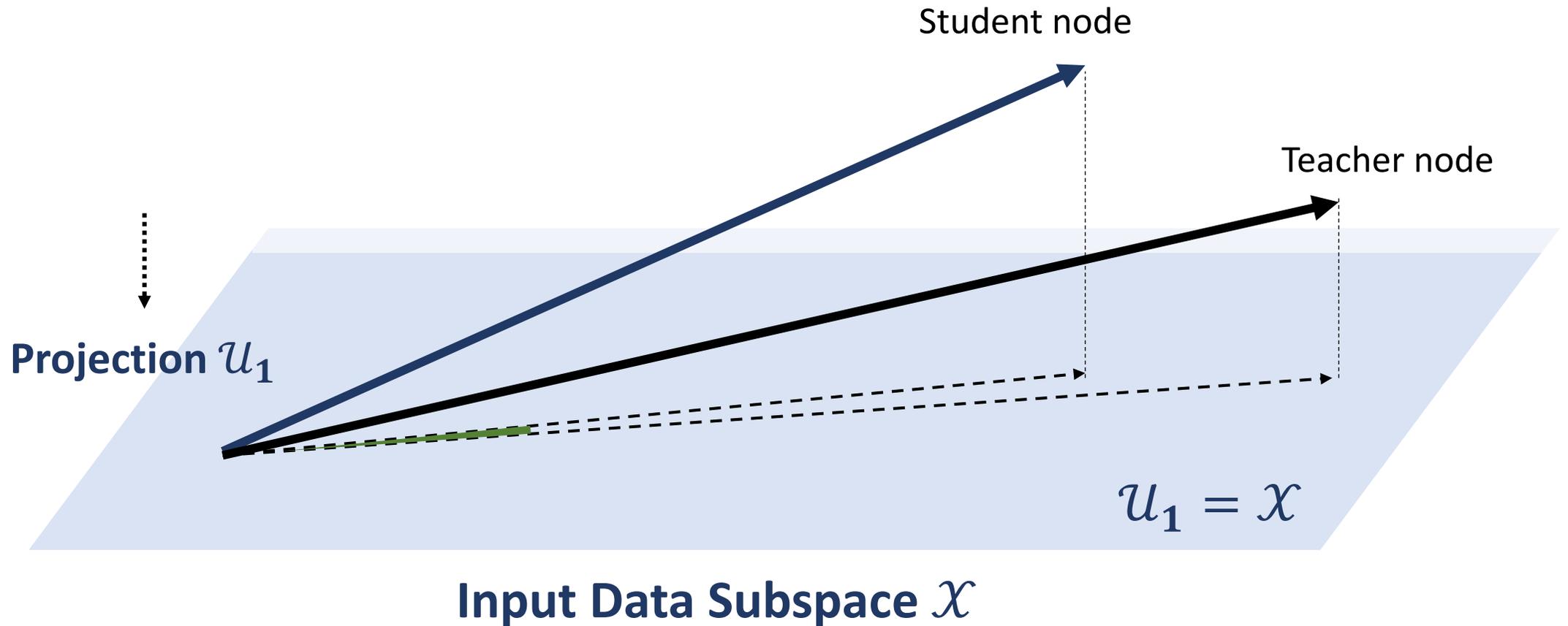
# Student Specialization

✓ **Full rank input space** → Specialization happens.  
Contributions of unspecialized nodes are zero.

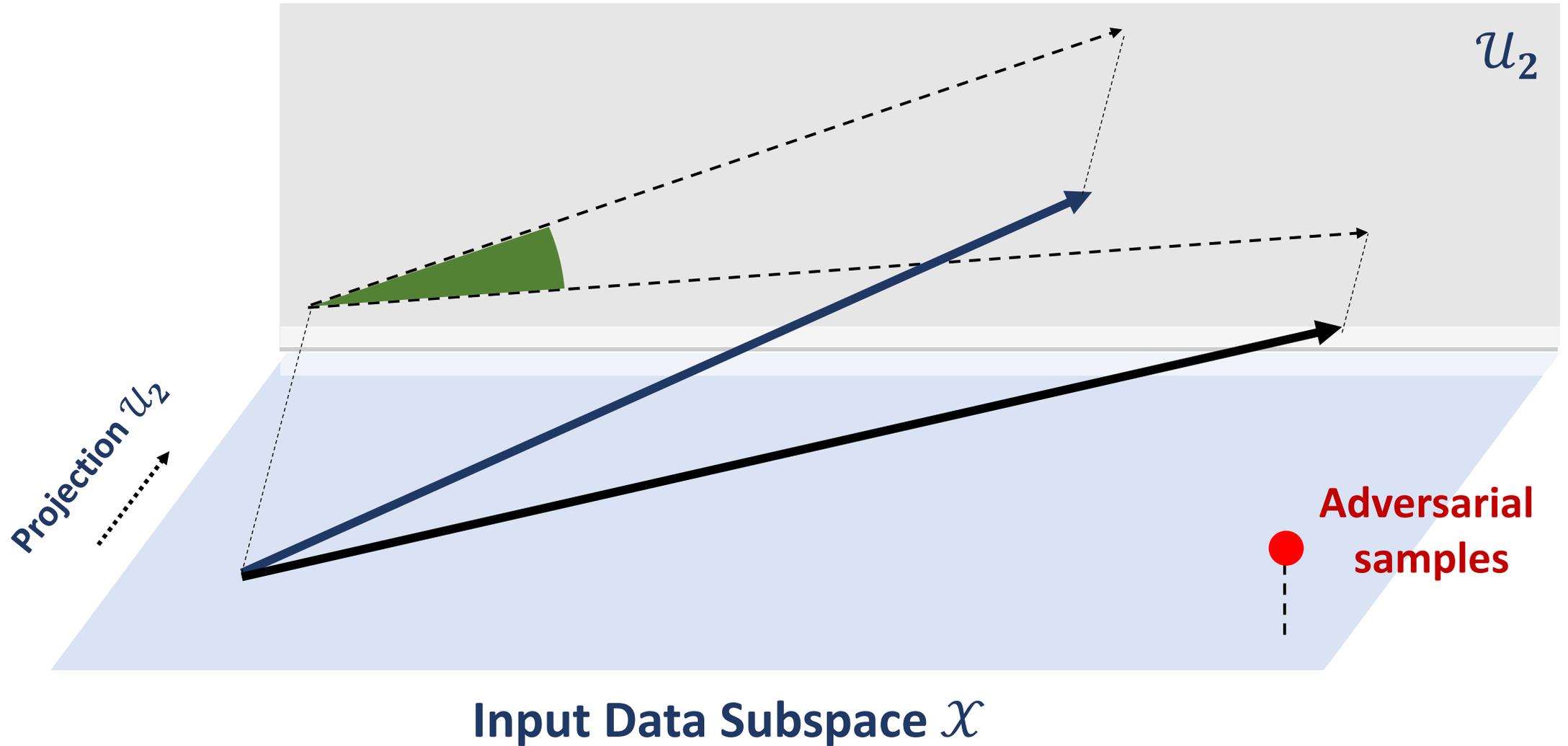
✗ **Low rank input** → ??



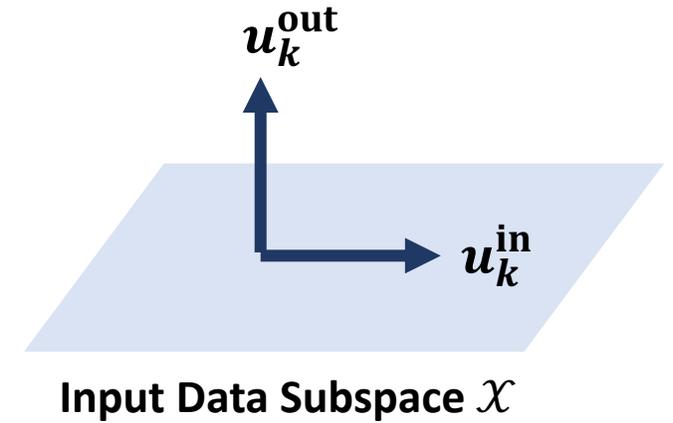
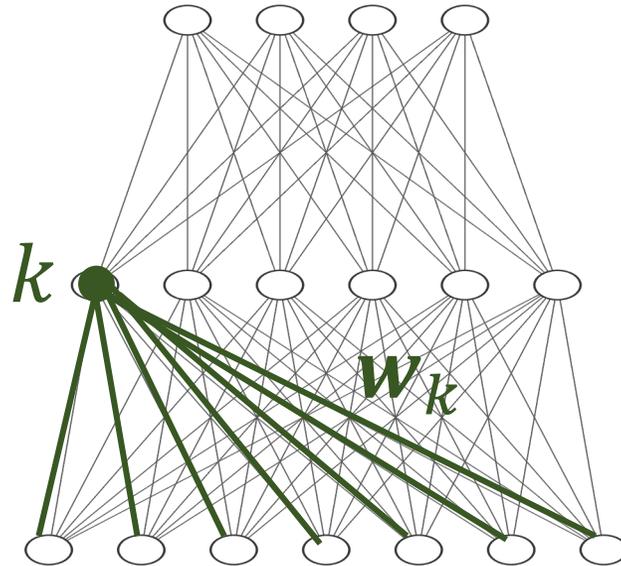
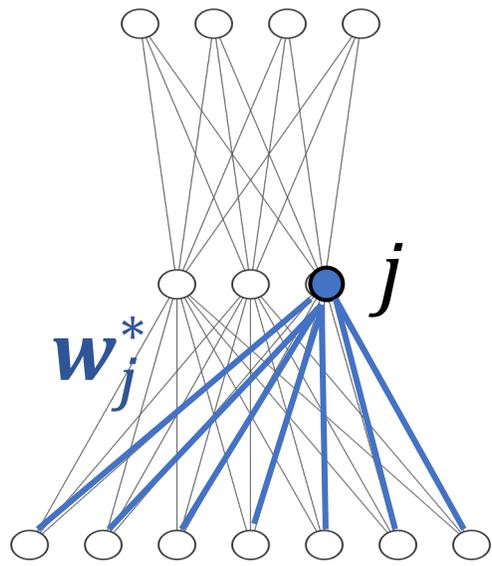
# Student Specialization under Low-rank Input



# Student Specialization under Low-rank Input



# An Empirical Model

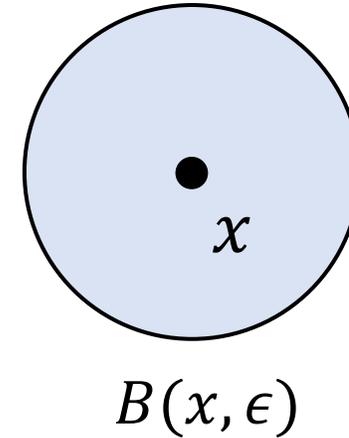
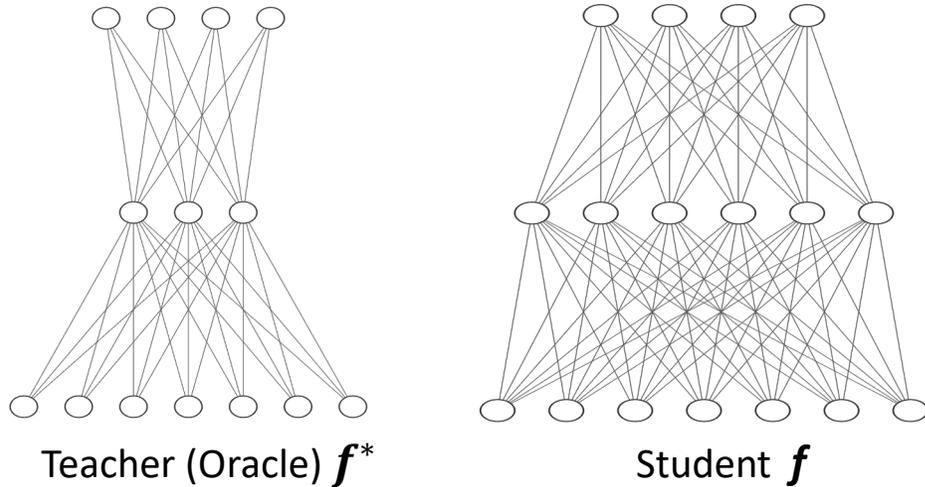


$$\mathbf{W}_k = \mathbf{w}_j^* + \epsilon_{\text{in}} \mathbf{u}_k^{\text{in}} + \epsilon_{\text{out}} \mathbf{u}_k^{\text{out}}$$

Student node k      Teacher node k

Will adversarial training improve the robust accuracy by remove the two?

# Revise Adversarial Samples



## Oracle Adversarial

$$\mathbf{x}' = \arg \max_{\mathbf{x}' \in B(\mathbf{x}, \epsilon)} L[f(\mathbf{x}'), f^*(\mathbf{x}')]$$

## Logit Training

$$\min_w L[f(x; w), f^*(x)]$$

## Data Adversarial

$$\mathbf{x}' = \arg \max_{\mathbf{x}' \in B(\mathbf{x}, \epsilon)} L[f(\mathbf{x}'), \mathbf{y}]$$

Fix target

## Label Training

$$\min_w L[f(x; w), \operatorname{argmax} f^*(x)]$$

# Experiments using CIFAR 10 images

Use CIFAR 10 images as low-dimensional and finite sample input

Use a 4-layered teacher network (trained on CIFAR10 and then pruned to be 45-32-32-20) as the *oracle*

## Using L2 loss and Oracle Adversarial

$$\mathbf{x}' = \arg \max_{\mathbf{x}' \in B(\mathbf{x}, \epsilon)} \|f(\mathbf{x}') - f^*(\mathbf{x}')\|^2$$

Train with pair  $(x', f^*(x'))$ . i.e., ground truth label  $f^*(x')$  for adversarial sample  $x'$

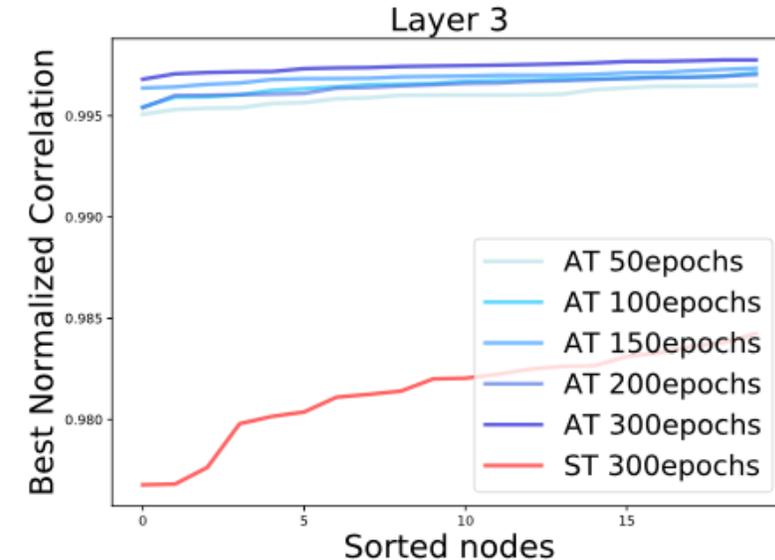
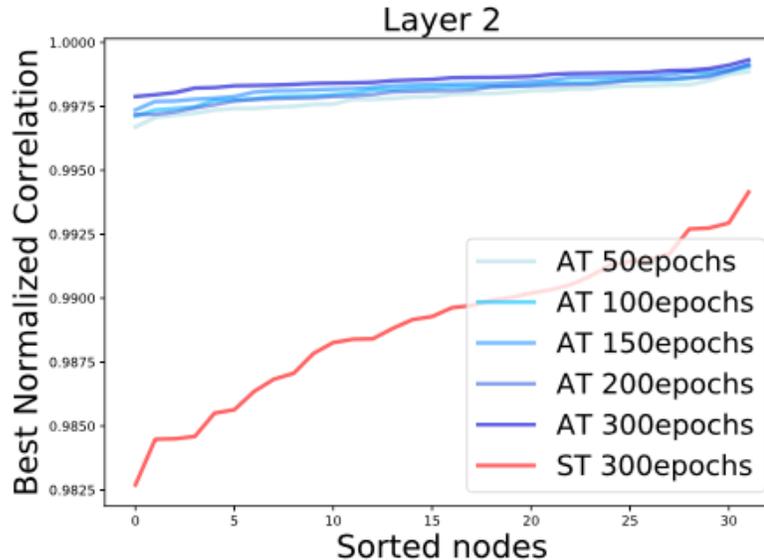
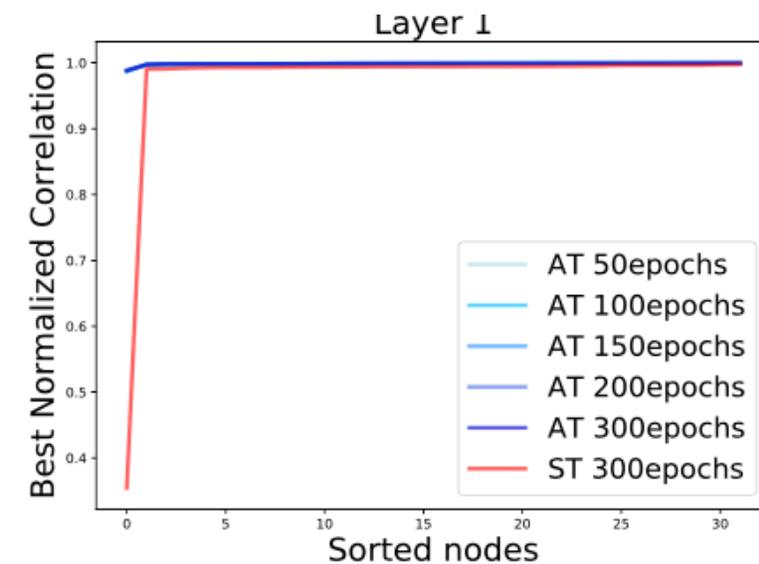
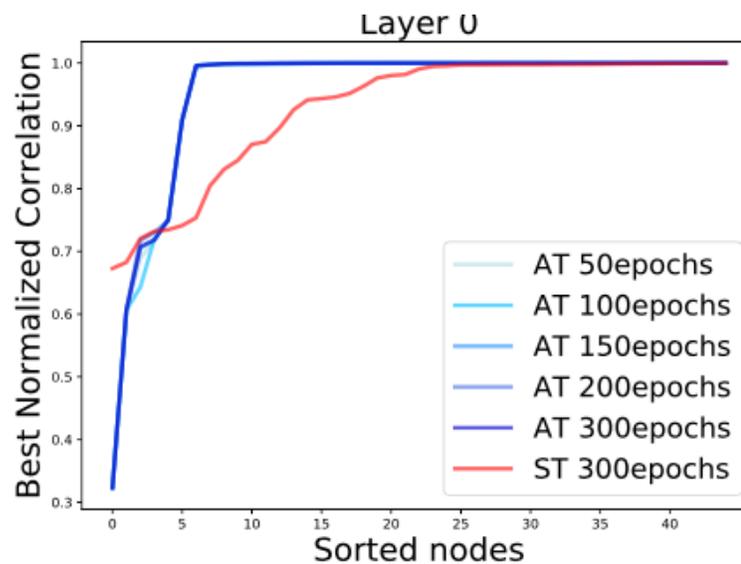
**(Use PGD in practice:** back-propagate  $f$  and  $f^*$  to get signed gradient, and apply to  $x$  iteratively)

# Experiments using CIFAR 10 images

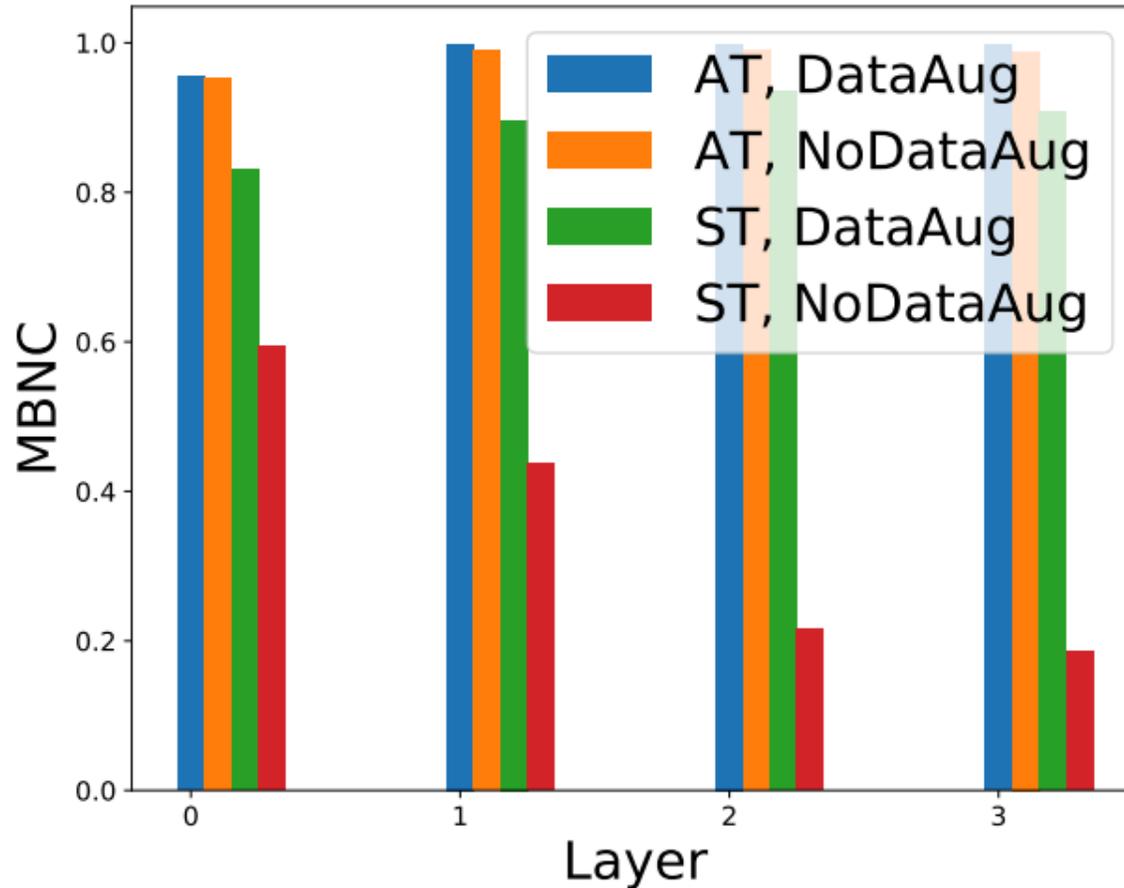
$\rho_{jk}$ : Normalized correlation between  $j$  and  $k$ :

$$\text{Best NC}(j) = \max_{\text{student } k} \rho_{jk}$$

Sorted Best NC curve

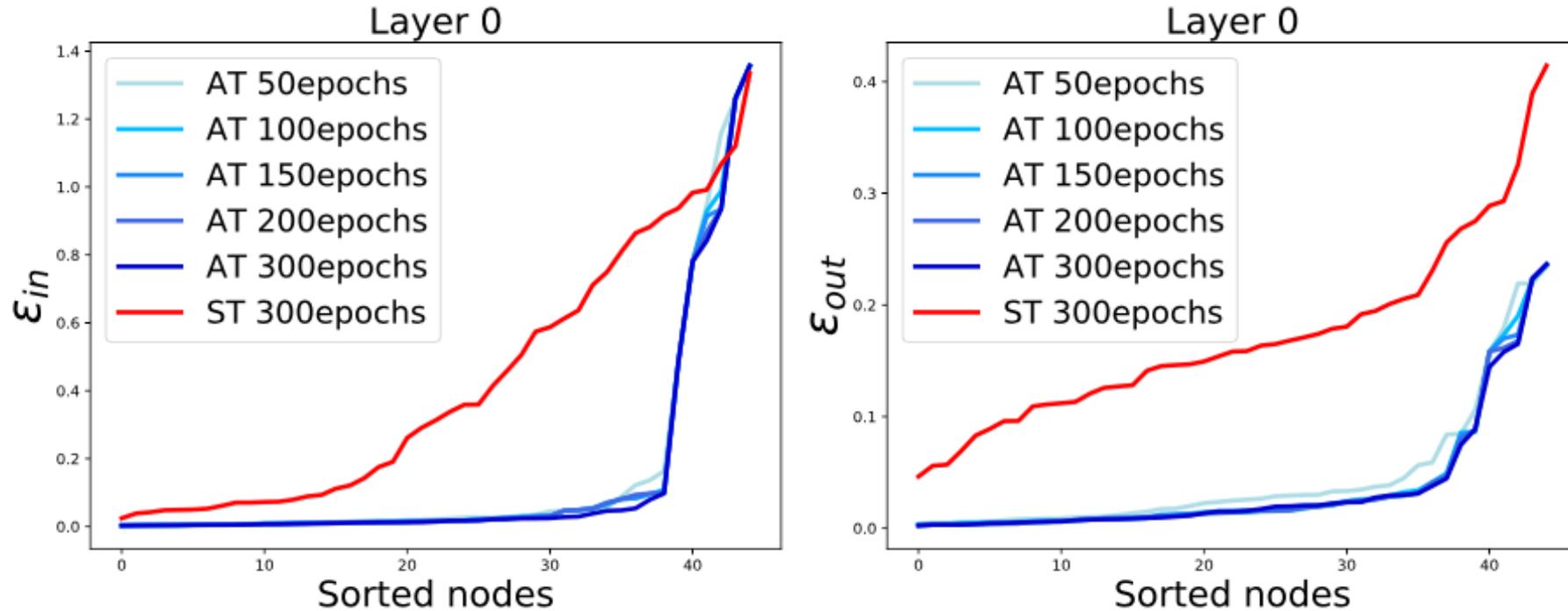


# Data Augmentation matters



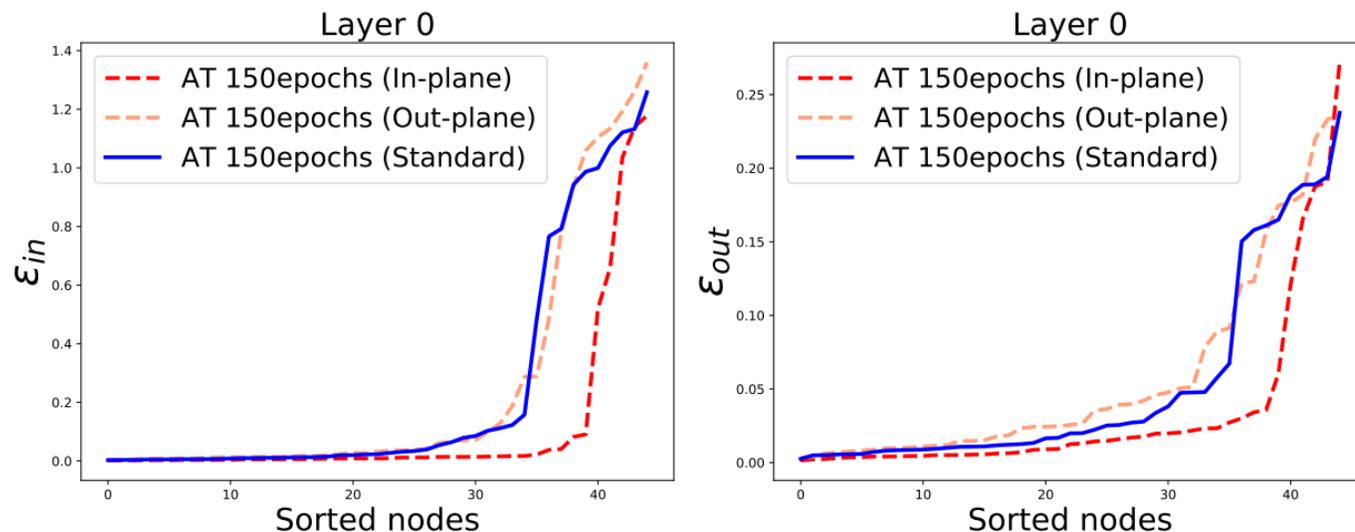
With data augmentation,  
we get better performance

# $\epsilon_{in}$ and $\epsilon_{out}$

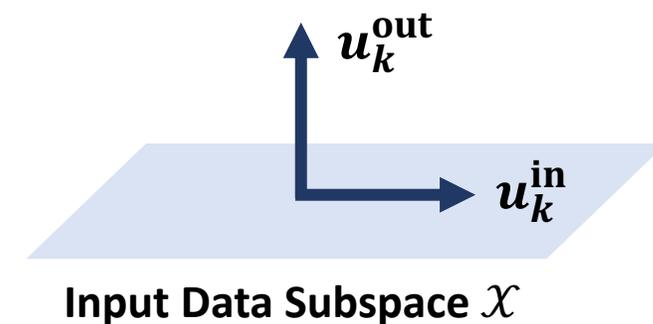


Model	AT (50)	AT (100)	AT (150)	AT (200)	AT (300)	ST (300)
Robust Acc	81.18%	82.57%	84.07%	85.35%	86.81%	62.77%

# Training with $\epsilon_{in}$ and $\epsilon_{out}$ only adversarial samples

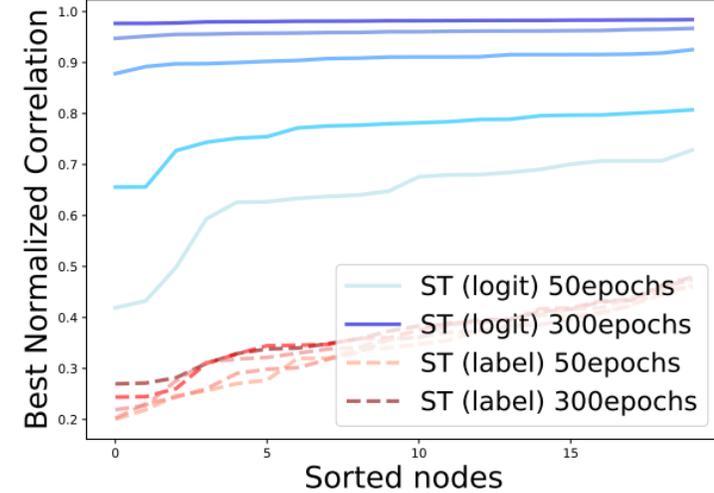
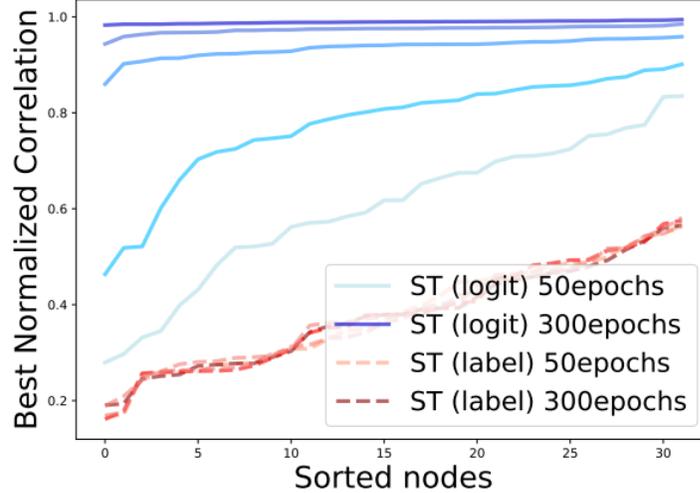
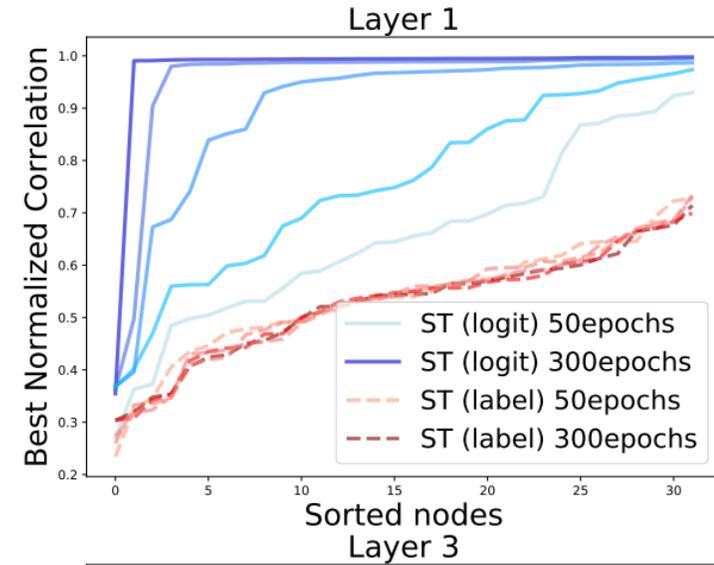
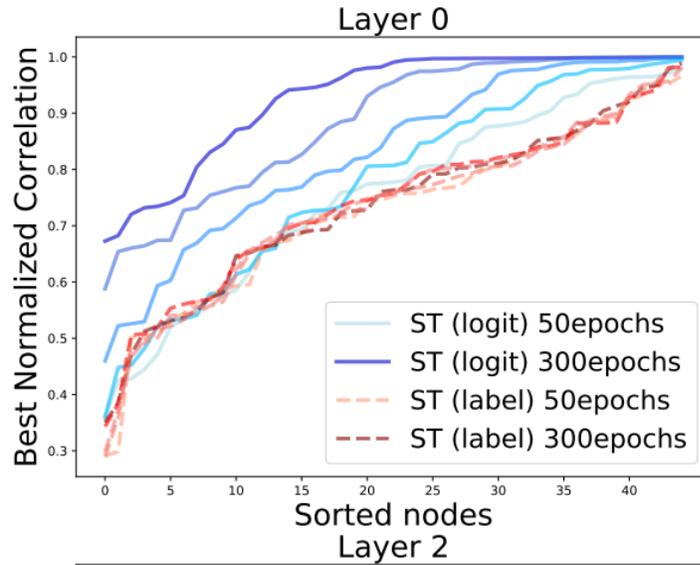


Attacks	In-plane	Out-plane	Standard
AT (In-plane)	88.86%	89.18%	89.28%
AT (Out-plane)	83.11%	83.54%	83.60%
AT (Standard)	86.87%	87.28%	87.18%

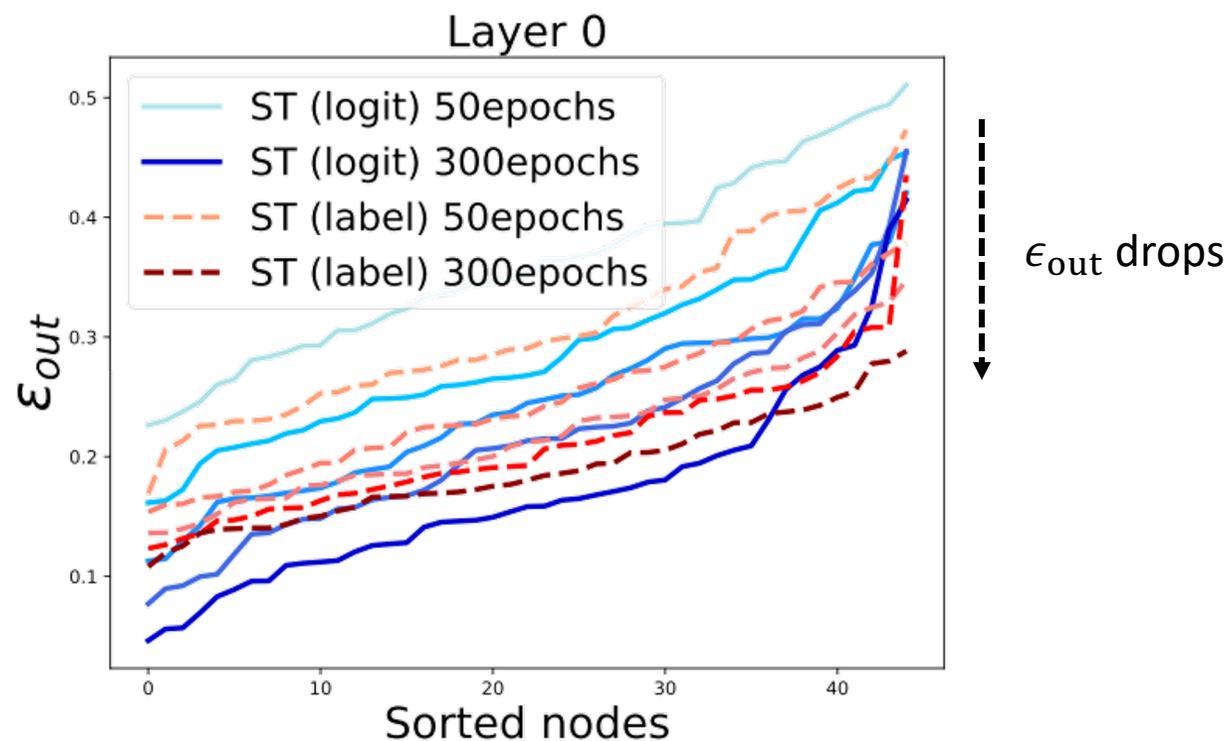
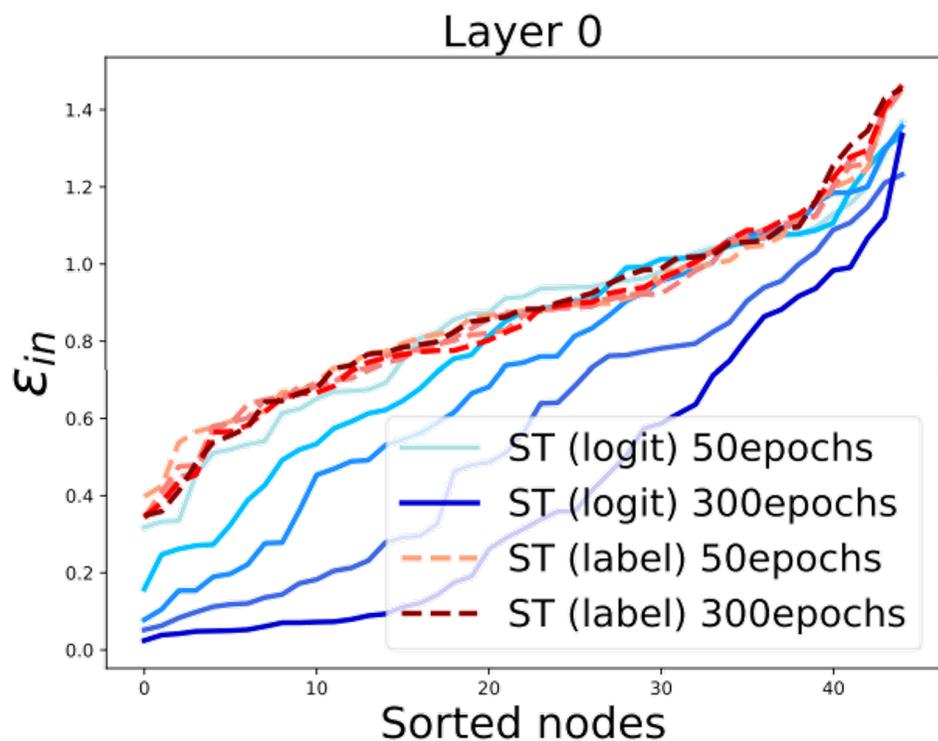


$\epsilon_{in}$  and  $\epsilon_{out}$  is only controlled in the first layer.

# Huge Difference between logits and label training



# Label training helps in $\epsilon_{out}$ but not $\epsilon_{in}$

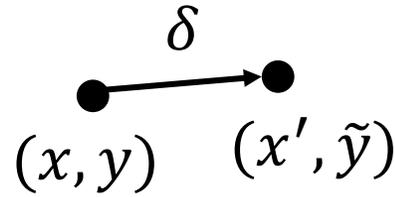


Robust Acc	ST (50)	ST (100)	ST (150)	ST (200)	ST (300)
Logit training	23.12%	30.72%	36.72%	48.52%	62.77%
Label training	19.08%	20.81%	22.34%	23.42%	25.79%

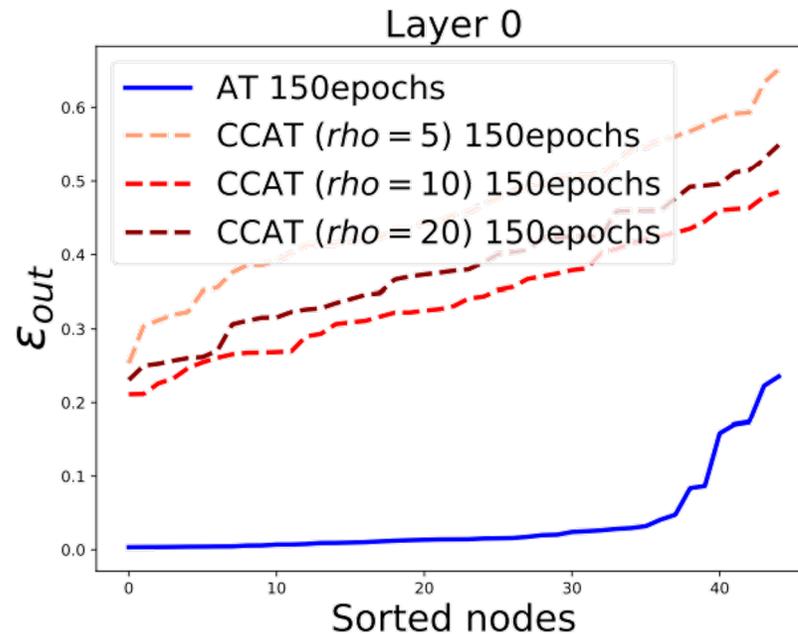
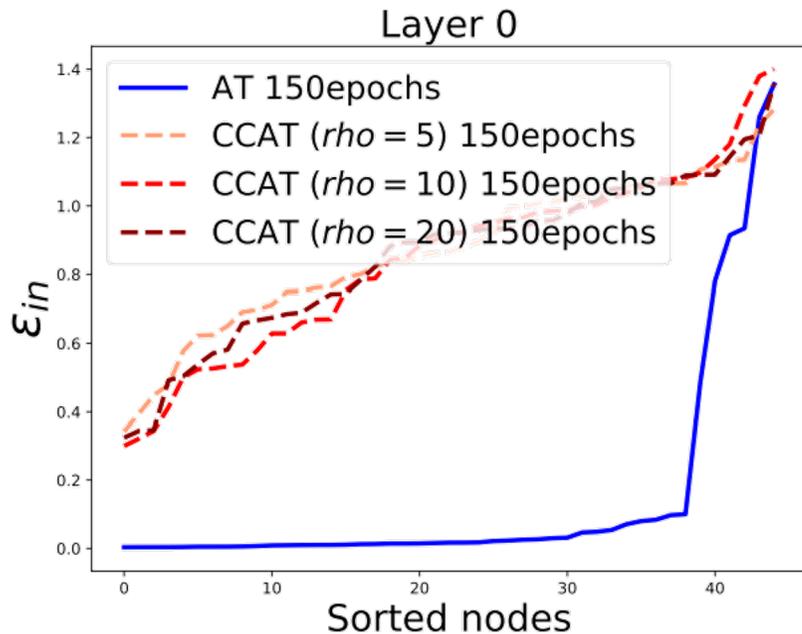
# Training with CCAT (AT without true label)

$$\tilde{y} = \lambda(\delta)\text{onehot}(y) + (1 - \lambda(\delta))\frac{1}{K}$$

$$\lambda(\delta) = (1 - \min(1, \|\delta\|_\infty/\epsilon))^\rho$$



Larger  $\delta$ , smaller  $\lambda$   
more uniform label  $\tilde{y}$



Model (150 epochs)	Robust Accuracy
CCAT ( $\rho = 5$ )	47.25%
CCAT ( $\rho = 10$ )	52.04%
CCAT ( $\rho = 20$ )	49.33%
Adversarial Training (logits)	84.07%

Thanks!