

Luck Matters: Understanding the Dynamics of Training Deep ReLU Neural Networks

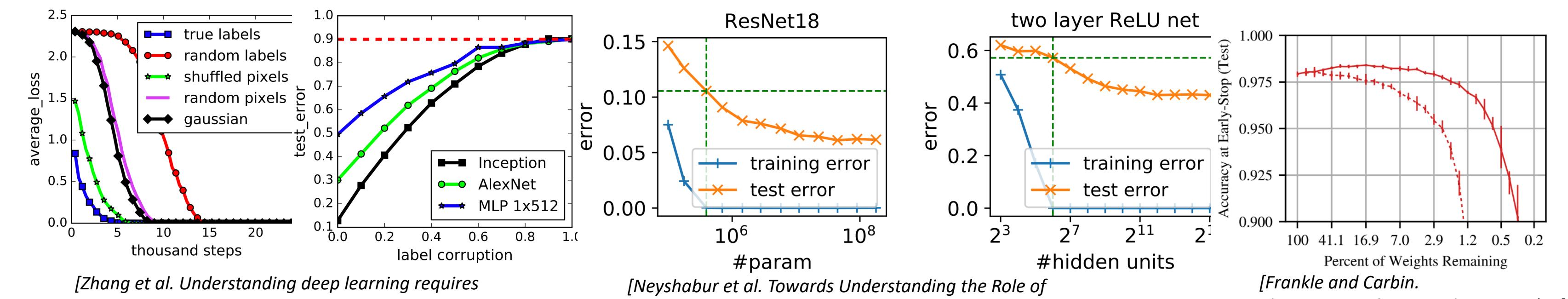
Yuandong Tian, Tina Jiang, Qucheng Gong, Ari Morcos

facebook

Artificial Intelligence Research

<https://github.com/facebookresearch/luckmatters>

Motivation



Over-Parameterization

More parameters, better test performance.
Network can be pruned substantially
Training with models with intrinsic capacity gives poor performance.

Implicit Regularization

Same network trained with SGD can fit both random and structured data
Network generalizes on structured data

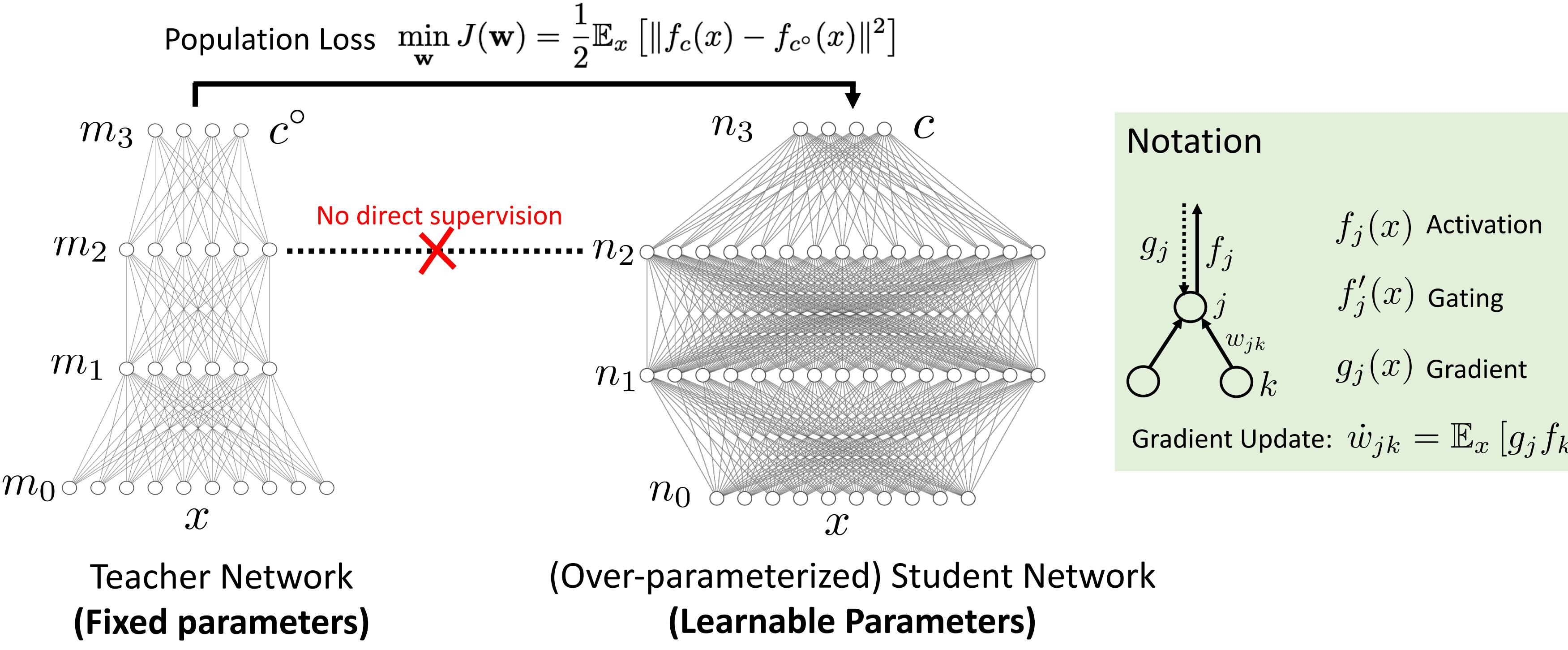
Lottery tickets

Use salient weights restarted to initialization gives lower test error
Use salient weights reinitialized gives poor performance.

Flat Minima

Many small eigenvalues in Hessian after convergence

Teacher-student setting



Recursive Gradient Rule

For the top-layer we have: $g_c(x) = f_{c^o}(x) - f_c(x)$

Is this condition apply to lower layers?

Base case:
 $\beta_{cc^o}^*(x) = \delta(c - c^o)$
 $\beta_{cc'}^*(x) = \delta(c - c')$

Theorem 1 Assuming for every node j in a layer, the gradient is:

$$g_j(x) = f'_j(x) \left[\sum_{j^o} \beta_{jj^o}^*(x) f_{j^o}(x) - \sum_{j'} \beta_{jj'}(x) f'_{j'}(x) \right] \beta_{kk^o}^*(x)$$

Compatibility between teacher k^o and student k

Then for the lower layer we have the same form with

$$\beta_{kk^o}^*(x) \equiv \sum_{jj^o} w_{jk} f'_j(x) \beta_{jj^o}^*(x) f'_{j^o}(x) w_{j^o k^o}^* \quad \beta_{kk'}(x) \equiv \sum_{jj'} w_{jk} f'_j(x) \beta_{jj'}(x) f'_{j'}(x) w_{j' k'}$$

Matrix Form of Gradient Descent

Student intermediate nodes mimics teacher

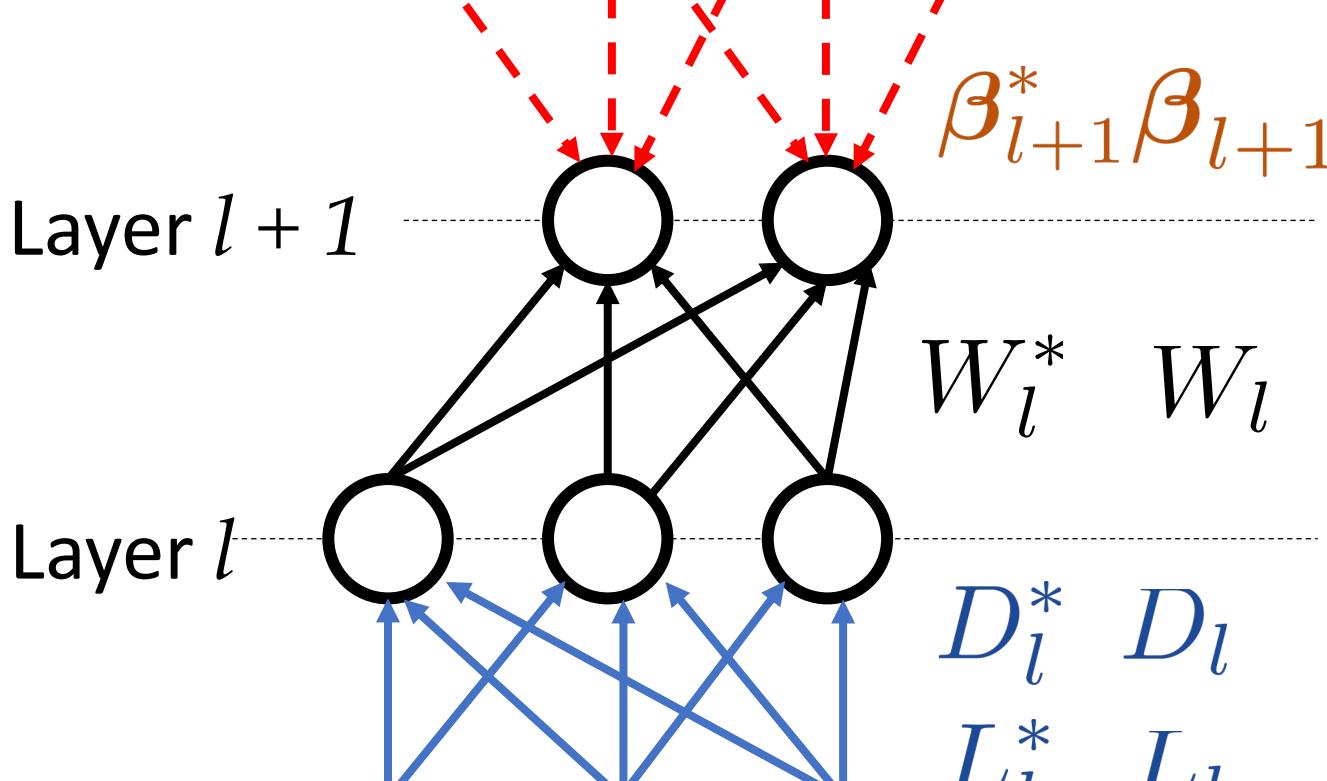
$$\dot{W}_l = L_l^* W_l^* H_{l+1}^* - L_l W_l H_{l+1}$$

$$[L^*]_{jj^o} = l_{jj^o}^* = E_x [f_j(x) f_{j^o}(x)] \quad [L]_{jj'} = l_{jj'} = E_x [f_j(x) f_{j'}(x)]$$

$$[D^*]_{jj^o} = d_{jj^o}^* = E_x [f'_j(x) f'_{j^o}(x)] \quad [D]_{jj'} = d_{jj'} = E_x [f'_j(x) f'_{j'}(x)]$$

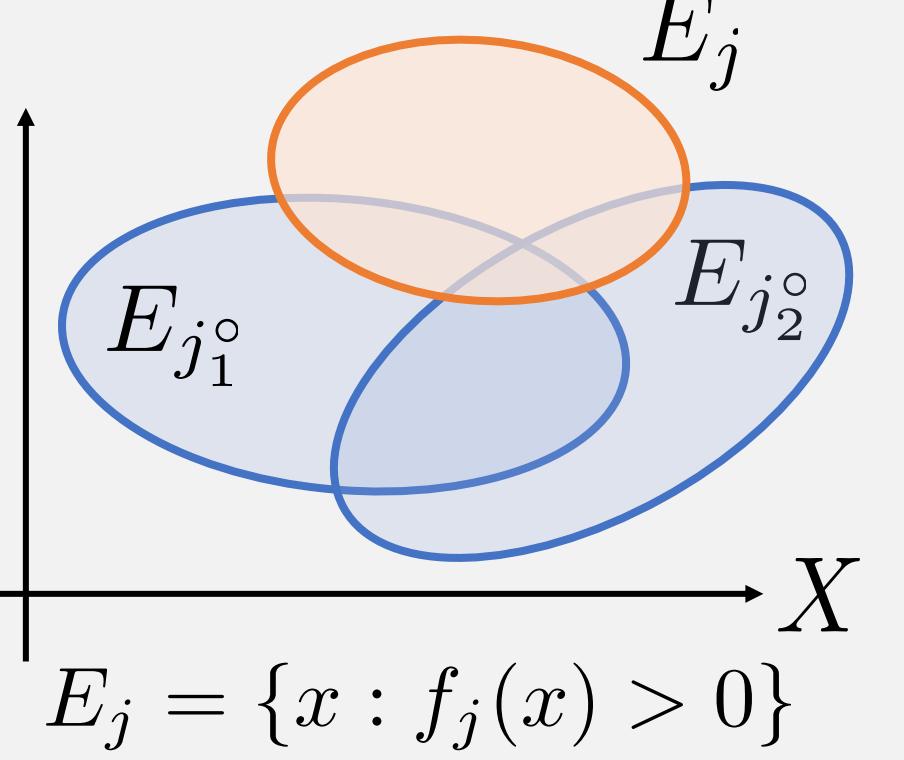
$$[\beta^*]_{jj^o} = E_x [\beta_{jj^o}^*(x)] \quad [\beta]_{jj'} = E_x [\beta_{jj'}(x)]$$

$$H^* = D^* \circ \beta^* \quad H = D \circ \beta$$

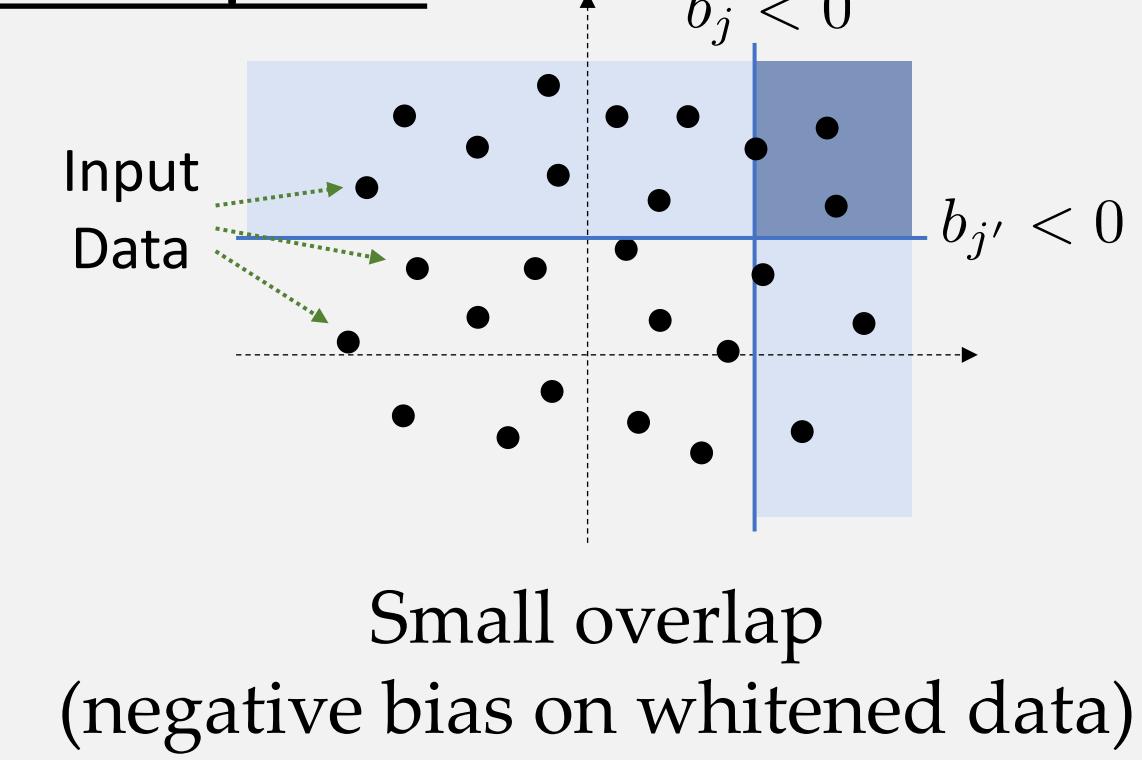


Main Result

Setting



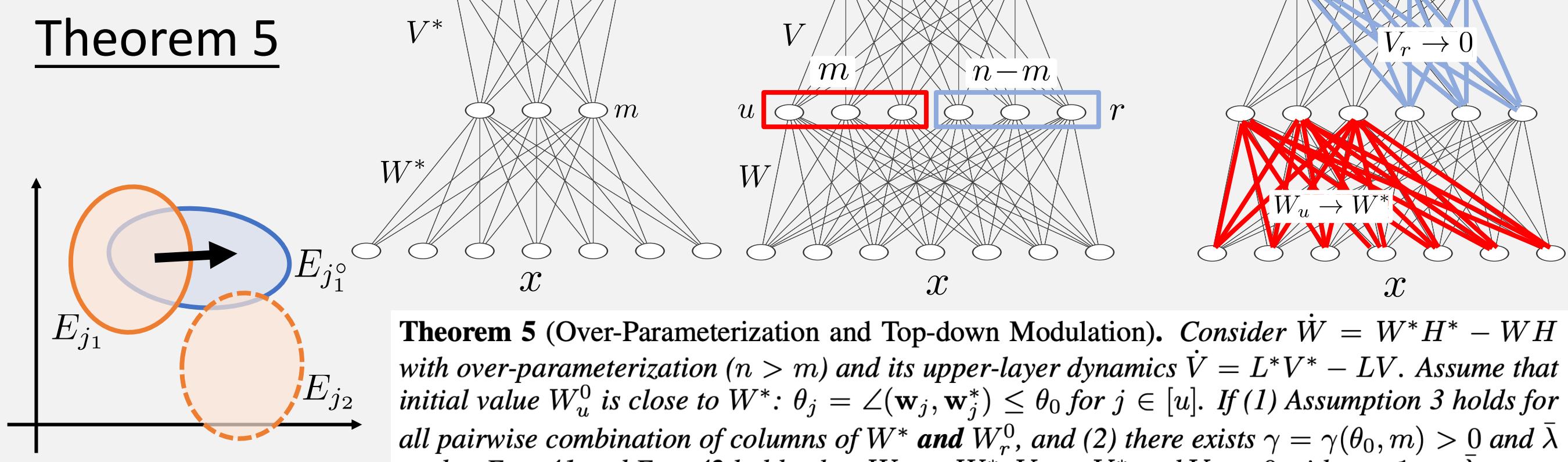
Assumption



Theorem 4

Theorem 4. For dynamics $\dot{w}_j = P_{w_j}^\perp (W^* h_j^* - Wh_j)$, where $P_{w_j}^\perp = I - w_j w_j^T$ is a projection matrix into the orthogonal complement of w_j . h_j^*, h_j are corresponding j -th column in H^* and H . Denote $\theta_j = \angle(w_j, w_j^*)$ and assume $\theta_j \leq \theta_0$. If $\gamma = \cos \theta_0 - (m-1)\epsilon M_d > 0$, then $w_j \rightarrow w_j^*$ with the rate $1 - \eta \bar{d} \gamma$ (η is learning rate). Here $\bar{d} = [1 + 2K_d \sin(\theta_0/2)] / \cos \theta_0$. $M_d = (1 + K_d)[1 + 2K_d \sin(\theta_0/2)]^2 / \cos^2 \theta_0$.

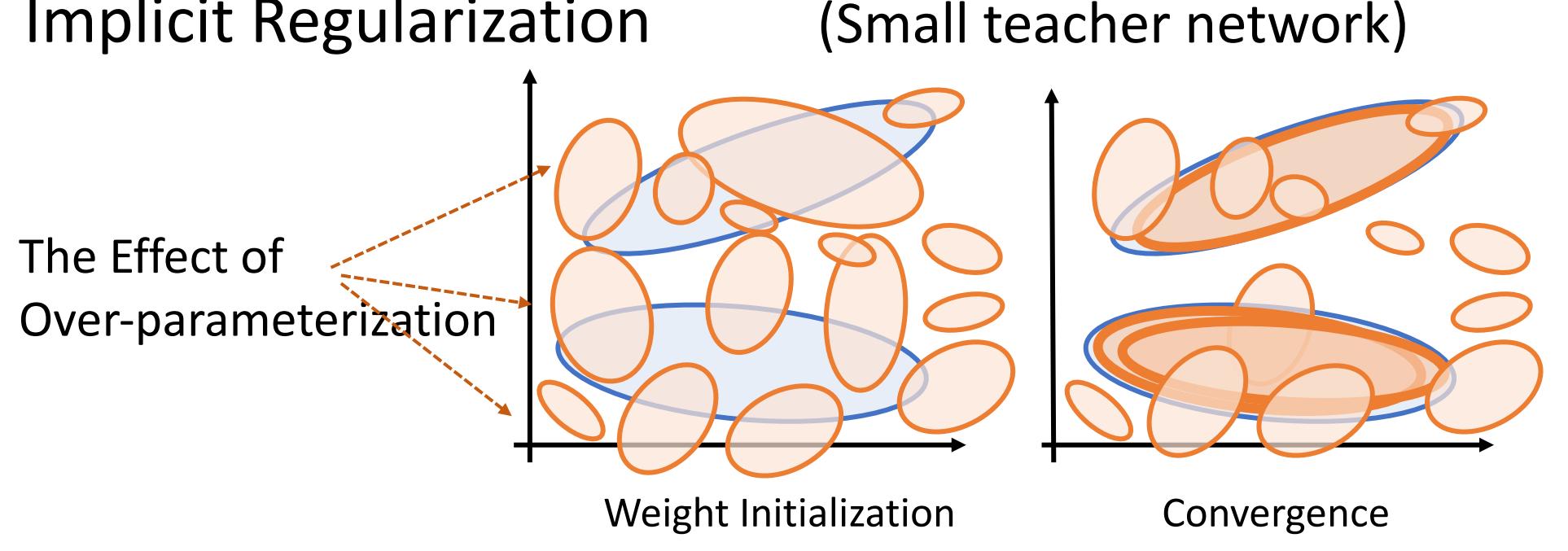
Theorem 5



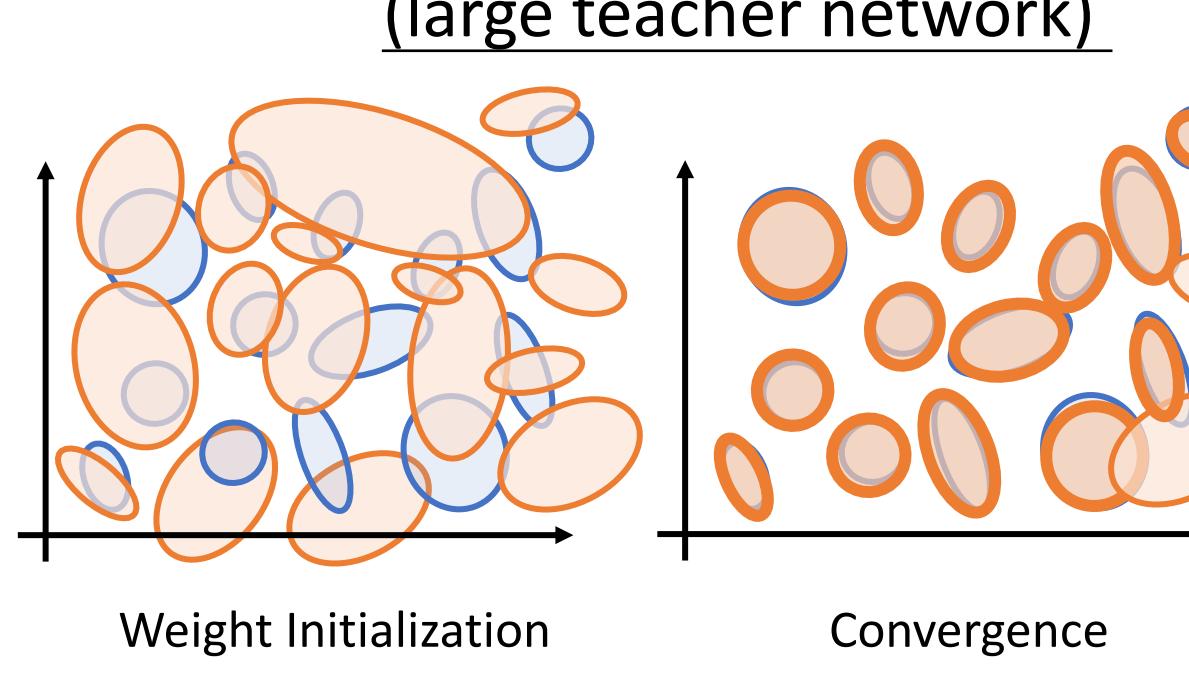
Theorem 5 (Over-Parameterization and Top-down Modulation). Consider $\dot{W} = W^* H^* - WH$ with over-parameterization ($n > m$) and its upper-layer dynamics $\dot{V} = L^* V^* - LV$. Assume that initial value W_u^0 is close to W^* : $\theta_j = \angle(w_j, w_j^*) \leq \theta_0$ for $j \in [u]$. If (1) Assumption 3 holds for all pairwise combination of columns of W^* and W_r^* , and (2) there exists $\gamma = \gamma(\theta_0, m) > 0$ and λ so that Eqn. 41 and Eqn. 42 holds, then $W_u \rightarrow W^*$, $V_u \rightarrow V^*$ and $V_r \rightarrow 0$ with rate $1 - \eta \bar{\lambda} \gamma$.

Explanation of Network Behaviors

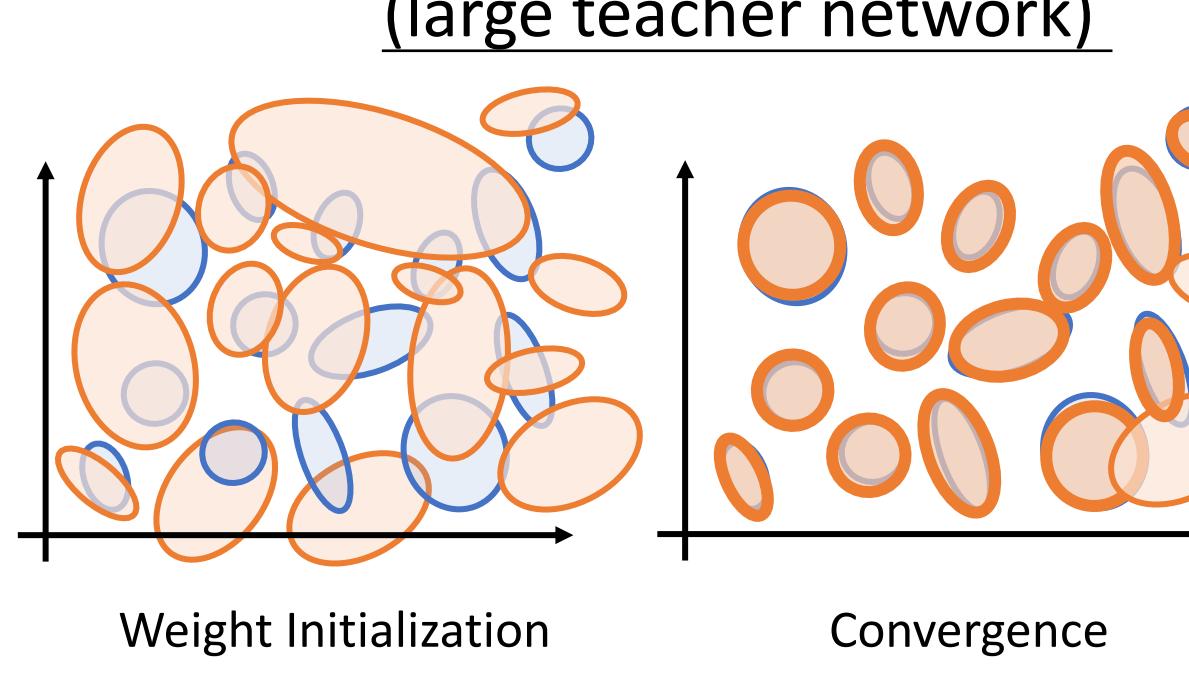
Implicit Regularization



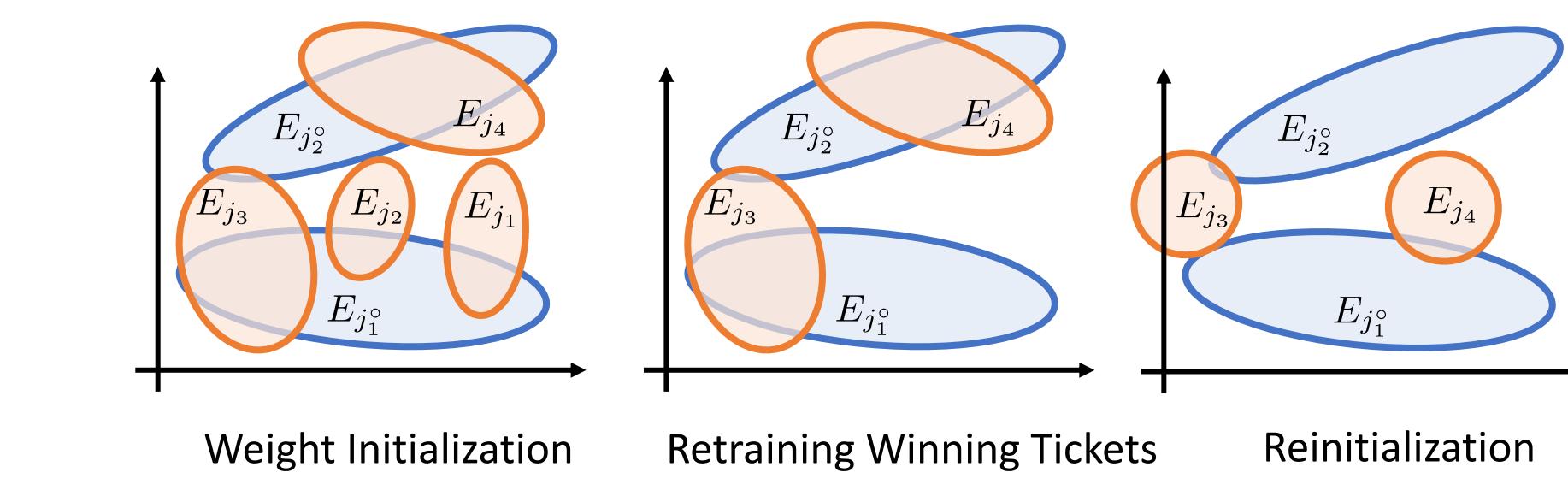
Structured Data (Small teacher network)



Random Data (large teacher network)



Lottery Tickets



Experiments

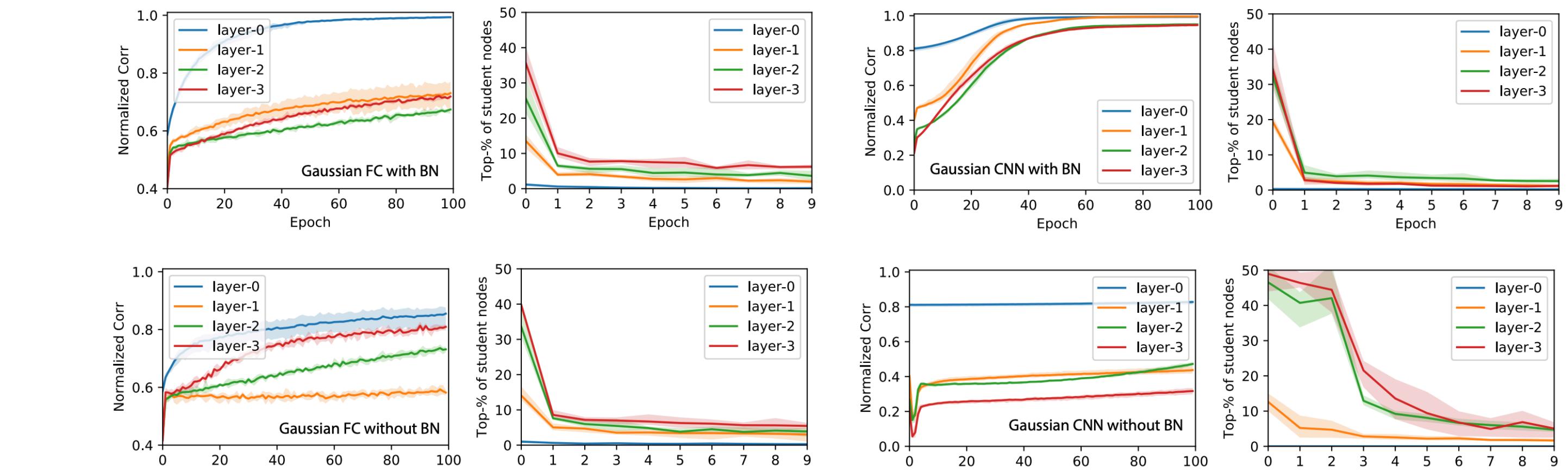
Setup

FC Teacher: 50-75-100-125
Convolution Teacher: 64-64-64-64
Student: 10x teacher

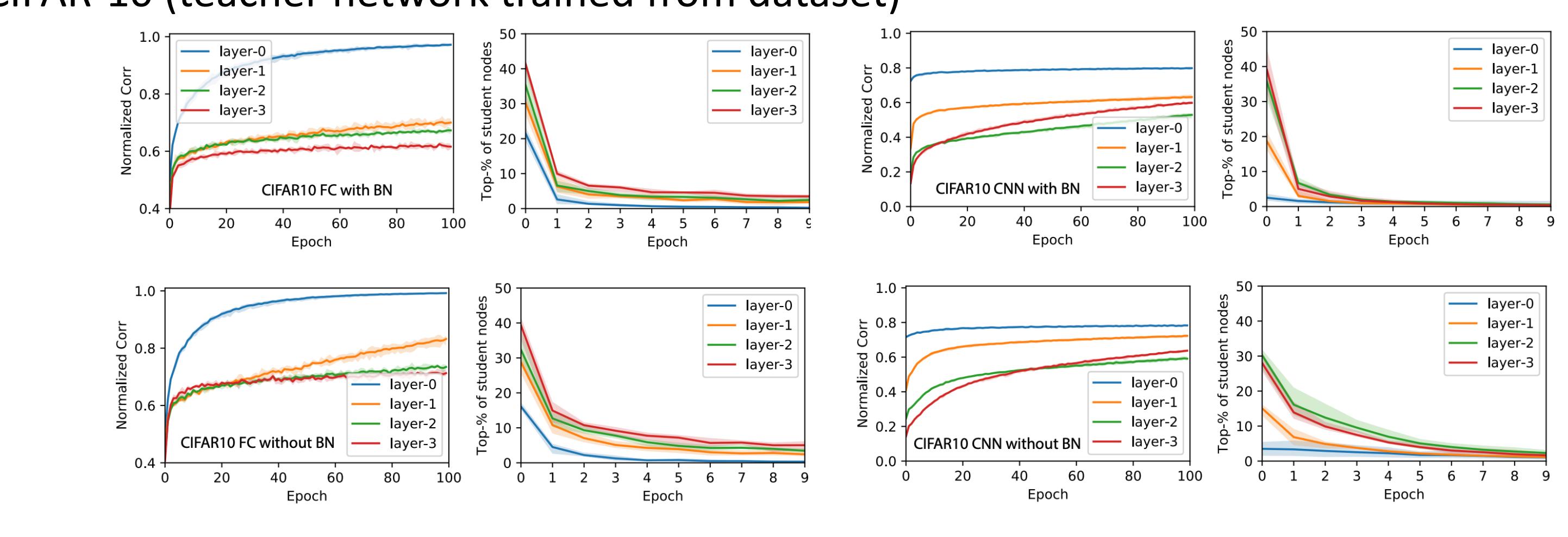
Performance Metric

Normalized Corr: $\bar{\rho} = \text{mean}_{j^o} \max_j f_j^T f_{j^o}$
Ranking of winner student node in early epochs

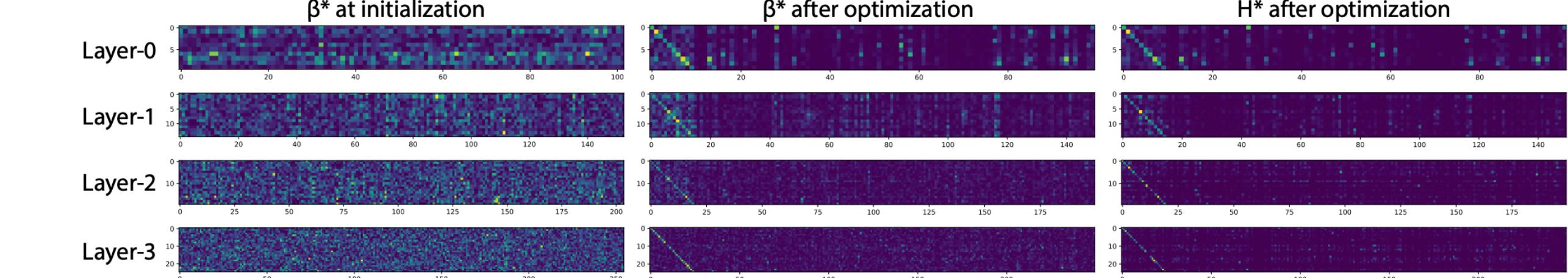
Random Gaussian Dataset: $x \sim N(0, 10I)$



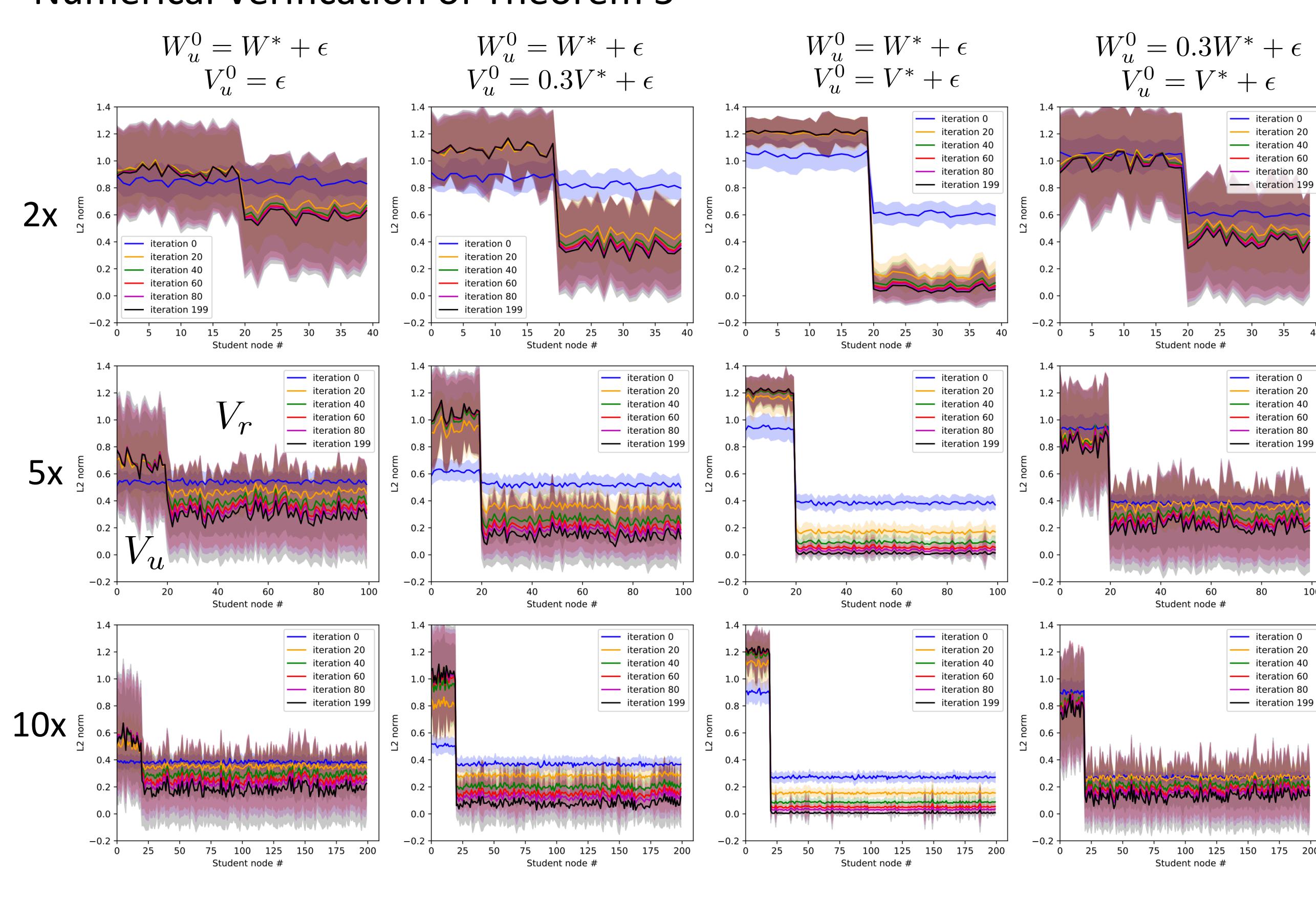
CIFAR-10 (teacher network trained from dataset)



Visualization



Numerical Verification of Theorem 5



Ablation Study

