

# Discussion on Neural Tangent Kernels (NTK) and Graph Neural Tangent Kernels (GNTK)

Prepared by Jiayao Zhang  
[jiayaozhang@acm.org](mailto:jiayaozhang@acm.org)

July 13, 2019

# Contents

## 1 Neural Tangent Kernels (NTK)

- Contruction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# Contents

## 1 Neural Tangent Kernels (NTK)

- Constrution of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# Overview

- The Neural Tangent Kernel (NTK) can be given by the kernel function

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{\theta}} \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle, \quad (1)$$

- where the gradient  $\frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}$  appears from the gradient descent.

$$f = \beta_0 + \beta_1 \cdot x$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\{\beta_0, \beta_1\}} [x^T x' + 1]$$

# Infinite Width Limit of the MLP (1/7)

- Write  $f(\boldsymbol{\theta}, \mathbf{x}) \in \mathbb{R}$  the output of a NN,  $\boldsymbol{\theta} \in \mathbb{R}^N$  all the parameters, and  $\mathbf{x} \in \mathbb{R}^d$  the input. Given a training dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ , consider training the neural network by minimizing the squared loss over training data:

$$\ell(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i)^2. \quad (2)$$

# Infinite Width Limit of the MLP (2/7)

- For an Multi-Layer Perceptron (MLP, feedforward nets with fully-connected layers) with  $L$  layers,  $d_1, \dots, d_L$  neurons per layer, and  $d_0 = d$ , where  $\mathbf{x} \in \mathbb{R}^d$ . Let  $\mathbf{x} \in \mathbb{R}^d$  be the input, write  $\mathbf{g}^{(0)}(\mathbf{x}) = \mathbf{x}$  and  $d_0 = d$ , and

$$\begin{aligned} \mathbf{f}^{(h)}(\mathbf{x}) &= \mathbf{W}^{(h)} \mathbf{g}^{(h-1)}(\mathbf{x}) \in \mathbb{R}^{d_h}, \\ \mathbf{g}^{(h)}(\mathbf{x}) &= \sqrt{\frac{c_\sigma}{d_h}} \sigma \left( \mathbf{f}^{(h)}(\mathbf{x}) \right) \in \mathbb{R}^{d_h}, \quad h = 1, 2, \dots, L. \end{aligned} \quad (3)$$

LeCun

where  $\mathbf{W}^{(h)} \in \mathbb{R}^{d_h \times d_{h-1}}$  is the weight matrix in the  $h$ -th layer ( $h \in [L]$ ),  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a coordinate-wise activation function, and

$$c_\sigma = \left( \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[ \sigma(z)^2 \right] \right)^{-1}. \quad \text{for ReLU, } c_{\sigma} = 2.$$

expected length / variance of all neurons.  
(Poole et al., 2016, exponential expressive)

# Infinite Width Limit of the MLP (3/7)

- The last layer of the neural network is

$$\begin{aligned}
 f(\boldsymbol{\theta}, \mathbf{x}) &= f^{(L+1)}(\mathbf{x}) = \mathbf{W}^{(L+1)} \cdot \mathbf{g}^{(L)}(\mathbf{x}) \\
 &= \mathbf{W}^{(L+1)} \cdot \sqrt{\frac{c_\sigma}{d_L}} \sigma \left( \mathbf{W}^{(L)} \cdot \sqrt{\frac{c_\sigma}{d_{L-1}}} \right. \\
 &\quad \left. \times \sigma \left( \mathbf{W}^{(L-1)} \dots \sqrt{\frac{c_\sigma}{d_1}} \sigma \left( \mathbf{W}^{(1)} \mathbf{x} \right) \right) \right), \tag{4}
 \end{aligned}$$

where  $\mathbf{W}^{(L+1)} \in \mathbb{R}^{1 \times d_L}$  is the weights in the final layer, and  $\boldsymbol{\theta} = (\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L+1)})$  represents all the parameters in the network.

# Infinite Width Limit of the MLP (4/7)

- All the weights are initialized as i.i.d.  $\mathcal{N}(0, 1)$ . In the limit of  $d_1, d_2, \dots, d_L \rightarrow \infty$ , the scaling factor  $\sqrt{c_\sigma/d_h}$  in Equation (4) ensures that the norm of  $\mathbf{g}^{(h)}(\mathbf{x})$  for each  $h \in [L]$  is approximately preserved at initialization (see <sup>#</sup>du2018global).
- In particular, for ReLU activation, we have  $\mathbb{E} \left[ \|\mathbf{g}^{(h)}(\mathbf{x})\|^2 \right] = \|\mathbf{x}\|^2$  ( $\forall h \in [L]$ ).

# Du et al. Gradient descent finds global minima of deep neural networks.  
<https://arxiv.org/abs/1811.03804>



# Infinite Width Limit of the MLP (5/7)

- From [LBN<sup>+</sup>17], one has the preactivations of each layer  $h \in [L]$  have their each coordinates tending to Gaussian process of covariance  $\Sigma^{(h-1)} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\begin{aligned}\Sigma^{(0)}(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^\top \mathbf{x}', \\ \mathbf{\Lambda}^{(h)}(\mathbf{x}, \mathbf{x}') &= \begin{pmatrix} \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \\ \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}') \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \\ \Sigma^{(h)}(\mathbf{x}, \mathbf{x}') &= c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\sigma(u) \sigma(v)],\end{aligned}\tag{5}$$

# Infinite Width Limit of the MLP (6/7)

- The intuition is that

$$\left[ \mathbf{f}^{(h+1)}(\mathbf{x}) \right]_i = \sum_{j=1}^{d_h} \left[ \mathbf{W}^{(h+1)} \right]_{i,j} \left[ \mathbf{g}^{(h)}(\mathbf{x}) \right]_j \quad (6)$$

is a centered Gaussian process conditioned on  $\mathbf{f}^{(h)}$  ( $\forall i \in [d_{h+1}]$ ), with covariance

$$\mathbb{E} \left[ \left[ \mathbf{f}^{(h+1)}(\mathbf{x}) \right]_i \cdot \left[ \mathbf{f}^{(h+1)}(\mathbf{x}') \right]_i \mid \mathbf{f}^{(h)} \right]. \quad (7)$$

# Infinite Width Limit of the MLP (7/7)

- We further have

$$\begin{aligned}
 & \text{if } i = k, \\
 & \text{Exp}[z^2] = 1; \\
 & \text{otherwise} \\
 & \text{Exp}[z_1 z_2] = 0; \\
 & z_1, z_2, z \sim \text{iid } \mathcal{N}(0, 1)
 \end{aligned}
 \quad
 \begin{aligned}
 & \mathbb{E} \left[ \left[ \mathbf{f}^{(h+1)}(\mathbf{x}) \right]_i \cdot \left[ \mathbf{f}^{(h+1)}(\mathbf{x}') \right]_i \middle| \mathbf{f}^{(h)} \right] \\
 &= \mathbb{E} \left[ \sum_{j, k \in [d_h]} \mathbf{W}_{ij}^{(h+1)} \mathbf{W}_{ik}^{(h+1)} \mathbf{g}^{(h)}(\mathbf{x})_j \mathbf{g}^{(h)}(\mathbf{x}')_k \middle| \mathbf{f}^{(h)} \right] \\
 &= \sum_{j, k \in [d_h]} \delta_{jk} \mathbf{g}^{(h)}(\mathbf{x})_j \mathbf{g}^{(h)}(\mathbf{x}')_k = \langle \mathbf{g}^{(h)}(\mathbf{x}), \mathbf{g}^{(h)}(\mathbf{x}') \rangle \\
 &= \frac{c_\sigma}{d_h} \langle \sigma(\mathbf{f}^{(h)}(\mathbf{x})), \sigma(\mathbf{f}^{(h)}(\mathbf{x}')) \rangle \\
 &= \frac{c_\sigma}{d_h} \sum_{j \in [d_h]} \sigma(\mathbf{f}^{(h)}(\mathbf{x}))_j \cdot \sigma(\mathbf{f}^{(h)}(\mathbf{x}'))_j \\
 &\xrightarrow[\text{a.s.}]{\text{LLN}} c_\sigma \mathbb{E}_{(u, v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\sigma(u) \sigma(v)] \equiv \Sigma^{(h)}(\mathbf{x}, \mathbf{x}'),
 \end{aligned} \tag{8}$$

as  $d_h \rightarrow \infty$  given that each  $\mathbf{f}_j^{(h)}$  is a centered Gaussian process with covariance  $\Sigma^{(h-1)}$ .

# Derivation of the NTK (1/7)

- To obtain the NTK, one computes the value that

$$\left\langle \frac{\partial f(\boldsymbol{\theta}, \boldsymbol{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \boldsymbol{x}')}{\partial \boldsymbol{\theta}} \right\rangle \quad (9)$$

converges to at random initialization in the infinite width limit.

# Derivation of the NTK (2/7)

- We can write the partial derivative with respect to a particular weight matrix  $\mathbf{W}^{(h)}$  in a compact form:

$$\frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \mathbf{W}^{(h)}} = \mathbf{b}^{(h)}(\mathbf{x}) \cdot \left( \mathbf{g}^{(h-1)}(\mathbf{x}) \right)^\top, \quad h = 1, 2, \dots, L+1, \quad (10)$$

(backprop)

where

$$\mathbf{b}^{(h)}(\mathbf{x}) = \begin{cases} 1 \in \mathbb{R}, & \text{for } h = L+1, \\ \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}) \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x}) \in \mathbb{R}^{d_h}, & \text{for } h \in [L]. \end{cases} \quad (11)$$

$$\mathbf{D}^{(h)}(\mathbf{x}) = \text{diag} \left( \dot{\sigma} \left( \mathbf{f}^{(h)}(\mathbf{x}) \right) \right) \in \mathbb{R}^{d_h \times d_h}, \quad h \in [L]. \quad (12)$$

# Derivation of the NTK (3/7)

- Then, for any  $h \in [L + 1]$ ,

$$\begin{aligned}
 \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \mathbf{W}^{(h)}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \mathbf{W}^{(h)}} \right\rangle &= \left\langle \mathbf{b}^{(h)}(\mathbf{x}) \cdot \left( \mathbf{g}^{(h-1)}(\mathbf{x}) \right)^\top, \right. \\
 &\quad \left. \mathbf{b}^{(h)}(\mathbf{x}') \cdot \left( \mathbf{g}^{(h-1)}(\mathbf{x}') \right)^\top \right\rangle \\
 &= \left\langle \mathbf{g}^{(h-1)}(\mathbf{x}), \mathbf{g}^{(h-1)}(\mathbf{x}') \right\rangle \\
 &\quad \times \left\langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \right\rangle.
 \end{aligned} \tag{13}$$

- Note that we have established in Equation (7) that

$$\left\langle \mathbf{g}^{(h-1)}(\mathbf{x}), \mathbf{g}^{(h-1)}(\mathbf{x}') \right\rangle \rightarrow \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}'). \tag{14}$$

# Derivation of the NTK (4/7)

- For the other factor  $\langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \rangle$ , by definition (11),

$$\begin{aligned} \langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \rangle &= \left\langle \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}) \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{h+1}(\mathbf{x}), \right. \\ &\quad \left. \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}') \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{h+1}(\mathbf{x}') \right\rangle. \end{aligned} \quad (15)$$

- In Eq. (11) of [ADH<sup>+</sup>19], there is another factor  $\frac{c_\sigma^{L-h}}{d_{h+1} \cdots d_L}$  preceeding the RHS of Eq. (15), which is likely to be a typo.

# Derivation of the NTK (5/7)

- Although  $\mathbf{W}^{(h+1)}$  and  $\mathbf{b}^{h+1}(\mathbf{x})$  are dependent, the Gaussian initialization of  $\mathbf{W}^{(h+1)}$  allows us to replace  $\mathbf{W}^{(h+1)}$  with a fresh new sample  $\widetilde{\mathbf{W}}^{(h+1)}$  without changing its limit.

$$\begin{aligned}
 & \left\langle \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}) (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{h+1}(\mathbf{x}), \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}') (\mathbf{W}^{(h+1)})^\top \mathbf{b}_{h+1}(\mathbf{x}') \right\rangle \\
 & \approx \left\langle \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}) (\widetilde{\mathbf{W}}^{(h+1)})^\top \mathbf{b}_{h+1}(\mathbf{x}), \sqrt{\frac{c_\sigma}{d_h}} \mathbf{D}^{(h)}(\mathbf{x}') (\widetilde{\mathbf{W}}^{(h+1)})^\top \mathbf{b}^{h+1}(\mathbf{x}') \right\rangle \\
 & \rightarrow \frac{c_\sigma}{d_h} \text{tr}(\mathbf{D}^{(h)}(\mathbf{x}) \mathbf{D}^{(h)}(\mathbf{x}')) \langle \mathbf{b}^{(h+1)}(\mathbf{x}), \mathbf{b}^{(h+1)}(\mathbf{x}') \rangle \\
 & \rightarrow \dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') \langle \mathbf{b}^{(h+1)}(\mathbf{x}), \mathbf{b}^{(h+1)}(\mathbf{x}') \rangle,
 \end{aligned} \tag{16}$$

where

$$\dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') = c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\dot{\sigma}(u) \dot{\sigma}(v)]. \tag{17}$$



# Derivation of the NTK (6/7)

- Applying this approximation inductively in Equation (15), we get

$$\left\langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \right\rangle \rightarrow \prod_{h'=h}^L \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}'). \quad (18)$$

- Finally, since

$$\left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle = \sum_{h=1}^{L+1} \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \mathbf{W}^{(h)}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \mathbf{W}^{(h)}} \right\rangle, \quad (19)$$

# Derivation of the NTK (7/7)

one has

$$\Theta^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{h=1}^{L+1} \left( \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}') \right), \quad (20)$$

NTK

where one writes  $\dot{\Sigma}^{(L+1)}(\mathbf{x}, \mathbf{x}') = 1$ .

# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# Convergence of NTK of MLP (1/2)

Theorem 3.1 [ADH<sup>+</sup>19]

Theorem 3.1 (Convergence to the NTK at initialization)

Fix  $\epsilon > 0$  and  $\delta \in (0, 1)$ . Suppose  $\sigma(z) = \max(0, z)$  ( $z \in \mathbb{R}$ ),  $\min_{h \in [L]} d_h \geq \text{poly}(L, 1/\epsilon) \cdot \log(L/\delta)$ , and  $[\mathbf{W}^{(h)}]_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$   $\forall h \in [L+1], i \in [d_h], j \in [d_{h-1}]$ . Then for any inputs  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$  such that  $\|\mathbf{x}\| \leq 1, \|\mathbf{x}'\| \leq 1$ , with probability at least  $1 - \delta$  we have:

$$\left| \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle - \Theta^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq \epsilon. \quad (21)$$

one realization of  
finite width

theoretical limit as  $d_h \rightarrow \infty$   
- dependent on data  
- indep of params

# Convergence of NTK of MLP (2/2)

## Theorem 3.1 [ADH<sup>+</sup>19]

- Compared with [JGH18] and [Yan19], Theorem 3.1 of [ADH<sup>+</sup>19] is non-asymptotic.
- [JGH18] requires  $d_0, \dots, d_L \rightarrow \infty$  sequentially; [Yan19] requires  $d_0, \dots, d_L \rightarrow \infty$  at the same rate. But [ADH<sup>+</sup>19] requires  $\min_h d_h \rightarrow \infty$ .

# Notations (1/8)

## Definition B.1 ( $k$ -homogeneous function)

A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is said to be  $k$ -homogeneous, if  $f(\lambda x) = \lambda^k f(x)$  for all  $x \in \mathbb{R}, \lambda > 0$ .

## Definition B.2

Let  $\mathcal{S}^+$  be the set of positive semi-definite kernels over  $\mathbb{R}^d$ , that is

$$\mathcal{S}^+ = \left\{ K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \middle| \forall N \in \mathbb{N}, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d, c_1, \dots, c_N \in \mathbb{R}, \right. \\ \left. \sum_{i=1}^N \sum_{j=1}^N c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0. \right\} \quad (22)$$

## Notations (2/8)

Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be the activation function, and  $\mathcal{T}_\sigma : \mathcal{S}^+ \rightarrow \mathcal{S}^+$  be the operator induced by  $\sigma$ , i.e.,

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d, \quad \mathcal{T}_\sigma(K)(\mathbf{x}, \mathbf{x}') = c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}|_{\mathbf{x}, \mathbf{x}'})} [\sigma(u) \sigma(v)], \quad (23)$$

where  $\mathbf{K}|_{\mathbf{x}, \mathbf{x}'} \in \mathbb{R}^{2 \times 2}$ ,  $\mathbf{K}|_{\mathbf{x}, \mathbf{x}'} = \begin{bmatrix} K(\mathbf{x}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}') \\ K(\mathbf{x}', \mathbf{x}) & K(\mathbf{x}', \mathbf{x}') \end{bmatrix}$ .

Define

$$t_\sigma(\Sigma) = c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \Sigma)} [\sigma(u) \sigma(v)], \quad (24)$$

$$\hat{t}_\sigma(\rho) \hat{t}_\sigma(\rho) = c_\sigma \mathbb{E}_{(u,v) \sim \Sigma'} [\sigma(u) \sigma(v)], \text{ with } \Sigma' = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \quad (25)$$

# Notations (3/8)

When  $\sigma$  is  $k$ -homogeneous function, we have

$$t_{\sigma}(\Sigma) = c_{\sigma} (\Sigma_{11}\Sigma_{22})^{\frac{k}{2}} \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \Sigma')} [\sigma(u) \sigma(v)], \quad (26)$$

with

$$\Sigma' = \begin{bmatrix} 1 & \frac{\Sigma_{12}}{\sqrt{\Sigma_{11}\Sigma_{22}}} \\ \frac{\Sigma_{12}}{\sqrt{\Sigma_{11}\Sigma_{22}}} & 1 \end{bmatrix}. \quad (27)$$

Thus  $t_{\sigma}(\Sigma)$  can be written as  $c_{\sigma} (\Sigma_{11}\Sigma_{22})^{\frac{k}{2}} \hat{t}(\frac{\Sigma_{12}}{\sqrt{\Sigma_{11}\Sigma_{22}}})$ ,



# Notations (4/8)

Fact B.1 (Some facts about  $\sigma(z) = \max(0, z)$  and  $\mathcal{T}_\sigma$ )

- ① For all activation function  $\sigma$ ,  $t_\sigma \left( \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right) = 1$ .
- ② For all  $\overset{\text{blue}}{k}$ -homogeneous activation  $\sigma$ ,  $\hat{t}_\sigma(1) = 1$  and  $t_\sigma \left( \begin{bmatrix} a & a \\ a & a \end{bmatrix} \right) = a^k$ .
- ③ For  $\sigma(z) = \max(0, z)$ ,  $\hat{t}_\sigma(\rho) = \frac{\sqrt{1-\rho^2} + \text{red}\arcsin \rho}{\pi} + \text{red}\frac{1}{2}$ ,  
 $\hat{t}_{\dot{\sigma}}(\rho) = \frac{1}{2} + \frac{\arcsin \rho}{\pi}$  and  $c_\sigma = c_{\dot{\sigma}} = 2$ .

(cf. 1905)

# Notations (5/8)

Recall the definition in Equation (5) and (17), we have

$$\begin{aligned}
 \Sigma^{(0)}(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^\top \mathbf{x}', \\
 \mathbf{\Lambda}^{(h)}(\mathbf{x}, \mathbf{x}') &= \Sigma^{(h-1)} \Big|_{\mathbf{x}, \mathbf{x}'} = \begin{pmatrix} \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \\ \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}') \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \\
 \Sigma^{(h)}(\mathbf{x}, \mathbf{x}') &= c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\sigma(u) \sigma(v)], \\
 \dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') &= c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\dot{\sigma}(u) \dot{\sigma}(v)]
 \end{aligned} \tag{28}$$

for  $h = 1, \dots, L$ .

For  $\sigma(z) = \max(z, 0)$ , we have

$$\Sigma^{(h)}(\mathbf{x}, \mathbf{x}) = \|\mathbf{x}\|^2, \quad \forall 0 \leq h \leq L. \tag{29}$$

# Notations (6/8)

Let  $\mathbf{D} = \mathbf{D}(\mathbf{x}, \mathbf{x}') = \mathbf{D}^{(h)}(\mathbf{x})\mathbf{D}^{(h)}(\mathbf{x}')$  is a 0-1 diagonal matrix. We define the following events:

- **Control of the post activations.**

- $\mathcal{A}^h(\mathbf{x}, \mathbf{x}', \epsilon_1) := \left\{ \left| \mathbf{g}^{(h)}(\mathbf{x}^{(0)})^\top \mathbf{g}^{(h)}(\mathbf{x}) - \Sigma^{(h)}(\mathbf{x}^{(0)}, \mathbf{x}) \right| \leq \epsilon_1 \right\},$   
 $\forall 0 \leq h \leq L$
- $\overline{\mathcal{A}}^h(\mathbf{x}, \mathbf{x}', \epsilon_1) = \mathcal{A}^h(\mathbf{x}, \mathbf{x}, \epsilon_1) \cap \mathcal{A}^h(\mathbf{x}, \mathbf{x}', \epsilon_1) \cap \mathcal{A}^h(\mathbf{x}', \mathbf{x}', \epsilon_1);$
- $\overline{\mathcal{A}}(\mathbf{x}, \mathbf{x}', \epsilon_1) = \bigcup_{h=0}^L \overline{\mathcal{A}}^h(\epsilon_1).$

- **Control of the backprop.**

- $\mathcal{B}^h(\mathbf{x}, \mathbf{x}', \epsilon_2) = \left\{ \left| \langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \rangle - \prod_{h=h}^L \dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') \right| < \epsilon_2 \right\};$
- $\overline{\mathcal{B}}^h(\mathbf{x}, \mathbf{x}', \epsilon_2) = \mathcal{B}^h(\mathbf{x}, \mathbf{x}, \epsilon_2) \cap \mathcal{B}^h(\mathbf{x}, \mathbf{x}', \epsilon_2) \cap \mathcal{B}^h(\mathbf{x}', \mathbf{x}', \epsilon_2);$
- $\overline{\mathcal{B}}(\mathbf{x}, \mathbf{x}', \epsilon_2) = \bigcup_{h=1}^{L+1} \overline{\mathcal{B}}^h(\epsilon_2);$

# Notations (7/8)

- **Control of the output.**

$$\overline{\mathcal{C}}(\mathbf{x}, \mathbf{x}', \epsilon_3) = \{|f(\boldsymbol{\theta}, \mathbf{x})| \leq \epsilon_3, |f(\boldsymbol{\theta}, \mathbf{x}')| \leq \epsilon_3\} \quad (30)$$

- **Control of the derivative.**

- $\mathcal{D}^h(\mathbf{x}, \mathbf{x}', \epsilon_4) = \left\{ \left| 2 \frac{\text{tr}(D(\mathbf{x}, \mathbf{x}'))}{d_h} - \dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') \right| < \epsilon_4 \right\};$
- $\overline{\mathcal{D}}^h(\mathbf{x}, \mathbf{x}', \epsilon_4) = \mathcal{D}^h(\mathbf{x}, \mathbf{x}, \epsilon_4) \cap \mathcal{D}^h(\mathbf{x}, \mathbf{x}', \epsilon_4) \cap \mathcal{D}^h(\mathbf{x}', \mathbf{x}', \epsilon_1);$
- $\overline{\mathcal{D}}(\mathbf{x}, \mathbf{x}', \epsilon_4) = \bigcup_{h=1}^{L+1} \overline{\mathcal{D}}^h(\epsilon_4).$

For simplicity, we will omit  $\mathbf{x}, \mathbf{x}'$  when there's no ambiguity. For events  $\mathcal{A}, \mathcal{B}$ , we define the event  $\mathcal{A} \Rightarrow \mathcal{B}$  as  $\neg \mathcal{A} \wedge \mathcal{B}$ .

# Notations (8/8)

Recall

$$\Theta^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{h=1}^{L+1} \left( \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}') \right), \quad (31)$$

where one writes  $\dot{\Sigma}^{(L+1)}(\mathbf{x}, \mathbf{x}') = 1$ ; and

## Theorem 3.1 (Convergence to the NTK at initialization)

Fix  $\epsilon > 0$  and  $\delta \in (0, 1)$ . Suppose  $\sigma(z) = \max(0, z)$  ( $z \in \mathbb{R}$ ),  $\min_{h \in [L]} d_h \geq \text{poly}(L, 1/\epsilon) \cdot \log(L/\delta)$ , and  $[\mathbf{W}^{(h)}]_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$   $\forall h \in [L+1], i \in [d_h], j \in [d_{h-1}]$ . Then for any inputs  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$  such that  $\|\mathbf{x}\| \leq 1, \|\mathbf{x}'\| \leq 1$ , with probability at least  $1 - \delta$  we have:

$$\left| \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle - \Theta^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq \epsilon. \quad (32)$$

# Proof Outline of Theorem 3.1 [ADH<sup>+</sup>19] (1/5)

## Theorem B.1 (Corollary 16 in [DFS16])

Let  $\sigma(z) = \max(0, z)$ ,  $z \in \mathbb{R}$  and  $[\mathbf{W}^{(h)}]_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ ,  
 $\forall h \in [L], i \in [d^{h+1}], j \in [d^h]$ , there exist constants  $c_1, c_2$ , such that if  
 $c_1 \frac{L^2 \log(\frac{8L}{\delta})}{\epsilon^2} \leq \min_{1 \leq h \leq L} d_h$  and  $\epsilon \leq \min(c_2, \frac{1}{L})$ , then for any fixed  
 $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ ,  $\|\mathbf{x}\|, \|\mathbf{x}'\| \leq 1$ , we have w.p.  $\geq 1 - \delta$ ,  $\forall 0 \leq h \leq L$ ,  
 $\forall (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,

$$\left| \mathbf{g}^{(h)}(\mathbf{x}^{(2)})^\top \mathbf{g}^{(h)}(\mathbf{x}^{(1)}) - \Sigma^{(h)}(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \right| \leq \epsilon. \quad (33)$$

In other words, if  $\min_{h \in [L]} d_h \geq c_1 \frac{L^2 \log(\frac{L}{\delta_1})}{\epsilon_1^2}$ ,  $\epsilon_1 \leq \min(c_2, \frac{1}{L})$ , then for  
fixed  $\mathbf{x}, \mathbf{x}'$ ,

$$\mathbb{P}[\text{there is at least one layer with the above ineq held}] \geq 1 - \delta_1.$$

# Proof Outline of Theorem 3.1 [ADH<sup>+</sup>19] (2/5)

## Theorem B.2

Let  $\sigma(z) = \max(0, z)$ ,  $z \in \mathbb{R}$ , if  $[\mathbf{W}^{(h)}]_{ij} \stackrel{i.i.d.}{\sim} N(0, 1)$ ,  
 $\forall h \in [L+1], i \in [d^{h+1}], j \in [d^h]$ , there exist constants  $c_1, c_2$ , such that if  
 $\min_{h \in [L]} d_h \geq c_1 \frac{L^2 \log(\frac{L}{\delta})}{\epsilon^4}$ ,  $\epsilon \leq \frac{c_2}{L}$ , then for any fixed  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ ,  
 $\|\mathbf{x}\|, \|\mathbf{x}'\| \leq 1$ , we have w.p.  $1 - \delta$ ,  $\forall 0 \leq h \leq L$ ,  
 $\forall (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,

$$\left| \mathbf{g}^{(h)}(\mathbf{x}^{(2)})^\top \mathbf{g}^{(h)}(\mathbf{x}^{(1)}) - \Sigma^{(h)}(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \right| \leq \frac{\epsilon^2}{2}, \quad (34)$$

and

$$\left| \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(1)}), \mathbf{b}^{(h)}(\mathbf{x}^{(2)}) \right\rangle - \prod_{h'=h}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| < 3L\epsilon. \quad (35)$$

# Proof Outline of Theorem 3.1 [ADH<sup>+</sup>19] (3/5)

## Theorem B.2

In other words, if  $\min_{h \in [L]} d_h \geq c_1 \frac{L^2 \log(\frac{L}{\delta_1})}{\epsilon_1^2}$ ,  $\epsilon_1 \leq \frac{c_2}{L}$ , then for fixed  $\mathbf{x}, \mathbf{x}'$ ,

$$\mathbb{P} \left[ \overline{\mathcal{A}} \left( \frac{\epsilon^2}{8} \right) \wedge \overline{\mathcal{B}}(3L\epsilon) \right] \geq 1 - \delta \quad (36)$$



# Proof Outline of Theorem 3.1 [ADH<sup>+</sup>19] (4/5)

## Proof of Theorem 3.1.

Recall that  $\Theta^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{h=1}^{L+1} \left( \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}') \right)$ , thus it suffices to show that w.p.  $1 - \epsilon$ , for every  $0 \leq h \leq L$ , it holds that

$$\left| \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \mathbf{W}^{(h)}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \mathbf{W}^{(h)}} \right\rangle - \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}') \right| \leq \frac{\epsilon}{L+1}, \quad (37)$$

which is a direct consequence of Theorem B.2 □

# Proof Outline of Theorem 3.1 [ADH<sup>+</sup>19] (5/5)

## Roadmap.

- Theorem 3.1  $\Leftarrow$  Theorem B.2;
- Theorem B.2  $\Leftarrow$  Theorem B.1, Lemmas B.4, B.6, B.7 (induction);
- Lemma B.7  $\Leftarrow$  Claims B.1, B.2, B.3;
- Claim B.2  $\Leftarrow$  Claim B.1, Lemmas B.3, B.10;
- Claim B.3  $\Leftarrow$  Lemma B.3;
- Lemma B.6  $\Leftarrow$  Lemma B.5  $\Leftarrow$  Lemmas B.8, B.9



# Some Lemmas (1/5)

## Lemma B.1 (Uniform Continuity of $\arcsin z$ )

- ① For any  $-\frac{\pi}{2} \leq y' \leq y \leq \frac{\pi}{2}$ ,  $\sin y - \sin y' \geq 2 \sin^2 \frac{y-y'}{2}$ .
- ②  $\sin y \geq \frac{2y}{\pi}, \forall y \in [0, \frac{\pi}{2}]$ .
- ③  $\arcsin$  is uniform continuous: for every  $\epsilon \in \mathbb{R}^+$ ,  
 $|z - z'| < \frac{2\epsilon^2}{\pi^2} \Rightarrow |\arcsin z - \arcsin z'| < \epsilon$ .
- ④ For  $\sigma(z) = \max(0, z)$ ,  $\hat{t}_{\dot{\sigma}}$  is uniform continuous: for every  $\epsilon \in \mathbb{R}^+$ ,  
 $|z - z'| < 2\epsilon^2 \Rightarrow |\hat{t}_{\dot{\sigma}}(z) - \hat{t}_{\dot{\sigma}}(z')| < \epsilon$ .

# Some Lemmas (2/5)

## Proof of Lemma B.1.

(1). From  $-\frac{\pi}{2} \leq y' \leq y \leq \frac{\pi}{2}$  we know  $\frac{-\pi}{2} + \frac{y-y'}{2} \leq \frac{y+y'}{2} \leq \frac{\pi}{2} - \frac{y-y'}{2}$ , which implies that  $\cos(\frac{y+y'}{2}) \geq \sin(\frac{y-y'}{2})$ . Thus,

$$\sin y \sin y' = 2 \cos \frac{y+y'}{2} \sin \frac{y-y'}{2} \geq 2 \sin^2 \frac{y-y'}{2}.$$

(2). Note that  $\left(\frac{\sin y}{y}\right)' = \frac{y \cos y - \sin y}{y^2} = \frac{\cos y}{y^2}(y - \tan y) < 0$ ,  $\frac{\sin y}{y}$  is decreasing on  $[0, \frac{\pi}{2}]$ . Thus  $\frac{\sin y}{y} \geq \frac{1}{\frac{\pi}{2}} = \frac{2}{\pi}, \forall y \in [0, \frac{\pi}{2}]$ .

(3). Let  $y, y' \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ , such that  $\sin y = z, \sin y' = z'$ . W.l.o.g., we assume  $y' < y, z' < z$ . Combing (1) and (2), we have

$$z - z' = \sin y - \sin y' \geq 2 \sin^2 \frac{y-y'}{2} \geq \frac{2(y-y')^2}{\pi^2}. \text{ Thus}$$

$$z - z' \leq \frac{2\epsilon^2}{\pi^2} \implies \arcsin z - \arcsin z' = y - y' \leq \epsilon.$$



# Some Lemmas (3/5)

## Lemma B.2

$$\mathbb{P}[\mathcal{A} \Rightarrow \mathcal{B}] \geq \mathbb{P}[\mathcal{B} \mid \mathcal{A}].$$

### Proof of Lemma B.2.

$$\begin{aligned} \mathbb{P}[\mathcal{A} \Rightarrow \mathcal{B}] &= \mathbb{P}[\neg \mathcal{A} \wedge \mathcal{B}] = 1 - \mathbb{P}[\mathcal{A} \vee \neg \mathcal{B}] = 1 - \mathbb{P}[\neg \mathcal{B} \mid \mathcal{A}] \mathbb{P}[\mathcal{A}] \geq \\ &1 - \mathbb{P}[\neg \mathcal{B} \mid \mathcal{A}] = \mathbb{P}[\mathcal{B} \mid \mathcal{A}]. \end{aligned}$$



For matrix  $\mathbf{A}$ , define the projection matrix for the column space of  $\mathbf{A}$ ,  $\mathbf{\Pi}_{\mathbf{A}} := \mathbf{A}\mathbf{A}^\dagger$  and the orthogonal projection matrix  $\mathbf{\Pi}_{\mathbf{A}}^\perp = I - \mathbf{A}\mathbf{A}^\dagger$ . For two random variables  $X$  and  $Y$ ,  $X \stackrel{\text{d}}{=}_{\mathcal{A}} Y$  means  $X$  is equal to  $Y$  in distribution conditioned on the  $\sigma$ -algebra generated by  $\mathcal{A}$ .

# Some Lemmas (4/5)

## Lemma B.3

Let  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_d)$ ,  $\mathbf{G} \in \mathbb{R}^{d \times k}$  be some fixed matrix, and random vector  $\mathbf{F} = \mathbf{w}^\top \mathbf{G}$ , then conditioned on the value of  $\mathbf{G}$ ,  $\mathbf{w}$  remains gaussian in the null space of the row space of  $\mathbf{G}$ . Mathematically, it means

$$\Pi_{\mathbf{G}}^\perp \mathbf{w} \stackrel{d}{=}_{\mathbf{F}=\mathbf{w}^\top \mathbf{G}} \Pi_{\mathbf{G}}^\perp \tilde{\mathbf{w}},$$

where  $\tilde{\mathbf{w}} \sim \mathcal{N}(0, \mathbf{I}_d)$  is a fresh i.i.d. copy of  $\mathbf{w}$ .

# Some Lemmas (5/5)

## Proof of Lemma B.3.

This lemma is straightforward when  $\Pi_G^\perp$  is a diagonal matrix. In general, let  $G = UG'$ , where  $U \in \mathbb{R}^{d \times d}$  is orthogonal and  $\Pi_{G'}^\perp$  is diagonal. Now we have

$$\Pi_G^\perp \mathbf{w} = U \Pi_{G'}^\perp U^\top \mathbf{w} \stackrel{\text{d}}{=}_{F=(U^\top \mathbf{w})^\top G'} U \Pi_{G'}^\perp U^\top \tilde{\mathbf{w}}, = \Pi_G^\perp \tilde{\mathbf{w}}$$

where we used the fact that if  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_d)$ , then for any orthogonal  $U$ ,  $U\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_d)$  twice. □

# Proof of Theorem B.2 (1/47)

## Theorem B.2

Let  $\sigma(z) = \max(0, z)$ ,  $z \in \mathbb{R}$ , if  $[\mathbf{W}^{(h)}]_{ij} \stackrel{i.i.d.}{\sim} N(0, 1)$ ,  
 $\forall h \in [L+1], i \in [d^{h+1}], j \in [d^h]$ , there exist constants  $c_1, c_2$ , such that if  
 $\min_{h \in [L]} d_h \geq c_1 \frac{L^2 \log(\frac{L}{\delta})}{\epsilon^4}$ ,  $\epsilon \leq \frac{c_2}{L}$ , then for any fixed  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ ,  
 $\|\mathbf{x}\|, \|\mathbf{x}'\| \leq 1$ , we have w.p.  $1 - \delta$ ,  $\forall 0 \leq h \leq L$ ,  
 $\forall (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,

$$\left| \mathbf{g}^{(h)}(\mathbf{x}^{(2)})^\top \mathbf{g}^{(h)}(\mathbf{x}^{(1)}) - \Sigma^{(h)}(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \right| \leq \frac{\epsilon^2}{2}, \quad (38)$$

and

$$\left| \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(1)}), \mathbf{b}^{(h)}(\mathbf{x}^{(2)}) \right\rangle - \prod_{h'=h}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| < 3L\epsilon. \quad (39)$$



# Proof of Theorem B.2 (2/47)

## Theorem B.2

*In other words, if  $\min_{h \in [L]} d_h \geq c_1 \frac{L^2 \log(\frac{L}{\delta_1})}{\epsilon_1^2}$ ,  $\epsilon_1 \leq \frac{c_2}{L}$ , then for fixed  $\mathbf{x}, \mathbf{x}'$ ,*

$$\mathbb{P} \left[ \overline{\mathcal{A}} \left( \frac{\epsilon^2}{8} \right) \wedge \overline{\mathcal{B}}(3L\epsilon) \right] \geq 1 - \delta \quad (40)$$

# Proof of Theorem B.2 (3/47)

Note that for  $c_\sigma = 2$  for  $\sigma(z) = \max(0, z)$ , by definition of  $\mathbf{b}^{(h)}$ , we have

$$\begin{aligned} & \langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \rangle \\ &= \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x})^\top \mathbf{W}^{(h+1)} \mathbf{D}^{(h)}(\mathbf{x}) \mathbf{D}^{(h)}(\mathbf{x}') \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x}'). \end{aligned} \quad (41)$$

Intuitively, when  $d_h$  is large, we can replace  $\mathbf{W}^{(h+1)}$  by a fresh i.i.d copy  $\widetilde{\mathbf{W}}$  with a small difference by  $\tilde{O}(\frac{1}{\sqrt{d_h}})$  as below. Similar techniques are used in [Yan19].

# Proof of Theorem B.2 (4/47)

$$\begin{aligned}
 \langle \mathbf{b}^{(h)}(\mathbf{x}), \mathbf{b}^{(h)}(\mathbf{x}') \rangle &= \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x})^\top \mathbf{W}^{(h+1)} \mathbf{D}^{(h)}(\mathbf{x}) \\
 &\quad \cdot \mathbf{D}^{(h)}(\mathbf{x}') (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}') \\
 &\approx \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x})^\top \widetilde{\mathbf{W}} \mathbf{D}^{(h)}(\mathbf{x}) \mathbf{D}^{(h)}(\mathbf{x}') \widetilde{\mathbf{W}}^\top \mathbf{b}^{(h+1)}(\mathbf{x}') \\
 &\approx \text{tr} \left( \frac{2}{d_h} \mathbf{D}^{(h)}(\mathbf{x}) \mathbf{D}^{(h)}(\mathbf{x}') \right) \mathbf{b}^{(h+1)}(\mathbf{x})^\top \mathbf{b}^{(h+1)}(\mathbf{x}') \\
 &\approx \dot{\Sigma}^{(h)}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \prod_{h'=h+1}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})
 \end{aligned} \tag{42}$$

The proof is based on a careful control of the following events.

# Proof of Theorem B.2 (5/47)

## Lemma B.4

$$\mathbb{P} \left[ \overline{\mathcal{A}}^L \left( \epsilon_1^2/2 \right) \implies \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{4}{\delta_3}} \right) \right] \geq 1 - \delta_3, \quad \forall \epsilon_1 \in [0, 1], \delta_3 \in (0, 1). \quad (43)$$

## Lemma B.5

$$\mathbb{P} \left[ \overline{\mathcal{A}}^{h+1} \left( \epsilon_1^2/2 \right) \implies \overline{\mathcal{D}}^h \left( \epsilon_1 + \sqrt{\frac{2 \log \frac{6}{\delta_4}}{d_h}} \right) \right] \geq 1 - \delta_4, \quad (44)$$

$$\forall \epsilon_1 \in [0, 1], \delta_4 \in (0, 1).$$

# Proof of Theorem B.2 (6/47)

## Lemma B.6

$$\mathbb{P} \left[ \overline{\mathcal{A}} \left( \epsilon_1^2/2 \right) \implies \overline{\mathcal{D}} \left( \epsilon_1 + \sqrt{\frac{2 \log \frac{6L}{\delta_4}}{\min_h d_h}} \right) \right] \geq 1 - \delta_4, \quad (45)$$

$$\forall \epsilon_1 \in [0, 1], \delta_4 \in (0, 1).$$

## Proof of Lemma B.6.

Apply union bound on Lemma B.5. □

# Proof of Theorem B.2 (7/47)

## Lemma B.7

There exists constant  $C, C' \in \mathbb{R}$ , for any  $\epsilon_2, \epsilon_3, \epsilon_4 \in [0, 1]$ , we have

$$\begin{aligned} & \mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon_1^2/2) \wedge \overline{\mathcal{B}}^{h+1} (\epsilon_2) \wedge \overline{\mathcal{C}} (\epsilon_3) \wedge \overline{\mathcal{D}}^h (\epsilon_4) \right. \\ & \quad \left. \implies \overline{\mathcal{B}}^h \left( \epsilon_2 + \frac{C' \epsilon_3}{\sqrt{d_h}} + 2\epsilon_4 + C \sqrt{\frac{\log \frac{1}{\delta_2}}{d_h}} \right) \right] \\ & \geq 1 - \delta_2. \end{aligned} \tag{46}$$

# Proof of Theorem B.2 (8/47)

## Proof of Theorem B.2.

We will use induction on Lemma B.7 to prove Theorem B.2. In the statement of Theorem B.1, we set  $\delta_1 = \frac{\delta}{4}$ ,  $\epsilon_1 = \frac{\epsilon^2}{8}$ , for some  $c_1, c_2$ , we have

$$\mathbb{P} \left[ \overline{\mathcal{A}}^L \left( \epsilon^2/8 \right) \right] \geq 1 - \delta/4 \quad (47)$$

# Proof of Theorem B.2 (9/47)

## Proof of Theorem B.2.

In the statement of Lemma B.6, we set  $\delta_4 = \frac{\delta_2}{4}$ , and  $\epsilon_1 = \frac{\epsilon}{2}$ . Note that for  $c_1$  large enough  $\sqrt{\frac{2 \log \frac{24L}{\delta}}{\min_h d_h}} \leq \frac{\epsilon}{2}$  and thus we have

$$\begin{aligned} \mathbb{P} \left[ \overline{\mathcal{A}} (\epsilon^2/8) \Rightarrow \overline{\mathcal{D}} (\epsilon) \right] &\geq \mathbb{P} \left[ \overline{\mathcal{A}} (\epsilon^2/8) \Rightarrow \overline{\mathcal{D}} \left( \epsilon/2 + \sqrt{\frac{2 \log \frac{24L}{\delta}}{\min_h d_h}} \right) \right] \\ &\geq 1 - \delta/4 \end{aligned} \quad (48)$$

In the statement of Lemma B.4, we set  $\delta_3 = \frac{\delta}{4}$ , and  $\epsilon_1 = \frac{\epsilon^2}{8}$ , we have

$$\mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon^2/8) \Rightarrow \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{16}{\delta}} \right) \right] \geq 1 - \delta/4 \quad (49)$$



# Proof of Theorem B.2 (10/47)

## Proof of Theorem B.2.

Using union bound on Equation (47),(48),(49), we have

$$\mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon^2/8) \wedge \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{16}{\delta}} \right) \wedge \overline{\mathcal{D}} (\epsilon) \right] \geq 1 - \frac{3\delta}{4} \quad (50)$$

Now we will begin the induction argument. First of all, note that  $\mathbb{P} \left[ \overline{\mathcal{B}}^{L+1} (0) \right] = 1$  by definition.

# Proof of Theorem B.2 (11/47)

## Proof of Theorem B.2.

For  $1 \leq h \leq L$  in the statement of Lemma B.7, we set  $\epsilon_2 = 3(L+1-h)\epsilon$ ,  $\epsilon_3 = 3\sqrt{\log \frac{16}{\delta}}$ ,  $\epsilon_2 = \epsilon$ ,  $\delta_4 = \frac{\delta}{4L}$ . Note that for  $c_1$  large enough,  $C\sqrt{\frac{\log \frac{1}{\delta_2}}{d_h}} + C'\sqrt{\frac{\log \frac{L}{\delta_2}}{d_h}} < \epsilon$ . Thus we have

# Proof of Theorem B.2 (12/47)

## Proof of Theorem B.2.

$$\begin{aligned}
 & \mathbb{P} \left[ \bar{\mathcal{B}}^{h+1} ((3L - 3h)\epsilon) \bigwedge \bar{\mathcal{C}} \left( 3\sqrt{\log \frac{16}{\delta_2}} \right) \bigwedge \bar{\mathcal{D}}^h (\epsilon) \right. \\
 & \quad \left. \Rightarrow \bar{\mathcal{B}}^h \left( (3L + 2 - 3h)\epsilon + C\sqrt{\frac{\log \frac{1}{\delta}}{d_h}} + 3C'\sqrt{\frac{\log \frac{16}{\delta}}{d_h}} \right) \right] \\
 & \geq \mathbb{P} \left[ \bar{\mathcal{B}}^{h+1} ((3L - 3h)\epsilon) \bigwedge \bar{\mathcal{C}} \left( 3\sqrt{\log \frac{16}{\delta_2}} \right) \bigwedge \bar{\mathcal{D}}^h (\epsilon) \right. \\
 & \quad \left. \Rightarrow \bar{\mathcal{B}}^h ((3L + 3 - 3h)\epsilon) \right] \\
 & \geq 1 - \frac{\delta}{4L}
 \end{aligned} \tag{51}$$

# Proof of Theorem B.2 (13/47)

## Proof of Theorem B.2.

Using union bound again on Equation (50) and Equation (51) for every  $h$  in  $\{1, 2, \dots, L\}$ , we have

$$\begin{aligned}
 & \mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon^2/8) \bigwedge \overline{\mathcal{C}} (\epsilon) \bigwedge \overline{\mathcal{D}} (\epsilon) \bigwedge \overline{\mathcal{B}} (3L\epsilon) \right] \\
 & \geq \mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon^2/8) \bigwedge \overline{\mathcal{C}} (\epsilon) \bigwedge \overline{\mathcal{D}} (\epsilon) \bigwedge_{h=1}^L \overline{\mathcal{B}} (3(L+1-h)\epsilon) \right] \\
 & \geq \left( 1 - \mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon^2/8) \bigwedge \overline{\mathcal{C}} (\epsilon) \bigwedge \overline{\mathcal{D}} (\epsilon) \right] \right) \\
 & \quad + \sum_{h=1}^L \left( 1 - \mathbb{P} \left[ \overline{\mathcal{B}}^{h+1} ((3L-3h)\epsilon) \bigwedge \overline{\mathcal{C}} (\epsilon) \bigwedge \overline{\mathcal{D}}^h (\epsilon) \Rightarrow \overline{\mathcal{B}}^h ((3L+3-3h)\epsilon) \right] \right) \\
 & \geq 1 - \delta
 \end{aligned} \tag{52}$$



# Proof of Theorem B.2 (14/47)

## Lemma B.4

$$\mathbb{P} \left[ \overline{\mathcal{A}}^L \left( \epsilon_1^2/2 \right) \implies \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{4}{\delta_3}} \right) \right] \geq 1 - \delta_3, \quad \forall \epsilon_1 \in [0, 1], \delta_3 \in (0, 1). \quad (53)$$

# Proof of Theorem B.2 (15/47)

## Proof of Lemma B.4.

For fixed  $\mathbf{g}^{(L)}(\mathbf{x})$ ,  $f(\boldsymbol{\theta}, \mathbf{x}) = \mathbf{W}^{(L+1)} \mathbf{g}^{(L)}(\mathbf{x}) \stackrel{d}{=} N(0, \|\mathbf{g}^{(L)}(\mathbf{x})\|^2)$ . Thus by subgaussian concentration[cite], we know w.p.  $\geq 1 - \delta$  over the randomness of  $\mathbf{W}^{(L+1)}$ ,  $|f(\boldsymbol{\theta}, \mathbf{x})| \leq \sqrt{2 \log \frac{2}{\delta}} \|\mathbf{g}^{(L)}(\mathbf{x})\|$ .

For  $\epsilon_1 \leq 1$ , we have  $\epsilon_1^2/2 < 1$ , which implies  $\|\mathbf{g}^{(L)}(\mathbf{x})\|^2 \leq 1 + \frac{\epsilon_1^2}{2} \leq 2$ , and thus taking union bound over  $\mathbf{x}, \mathbf{x}'$ , we have w.p.  $\geq 1 - \delta$ ,  
 $|f(\boldsymbol{\theta}, \mathbf{x})| \leq 2\sqrt{\log \frac{2}{\delta}}, |f(\boldsymbol{\theta}, \mathbf{x}')| \leq 2\sqrt{\log \frac{2}{\delta}}.$

# Proof of Theorem B.2 (16/47)

## Proof of Lemma B.4.

$$\begin{aligned} \mathbb{P} \left[ \overline{\mathcal{A}}^L \left( \epsilon_1^2/2 \right) \Rightarrow \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{4}{\delta_3}} \right) \right] &\geq \mathbb{P} \left[ \overline{\mathcal{C}} \left( 2\sqrt{\log \frac{4}{\delta_3}} \right) \mid \overline{\mathcal{A}}^L \left( \epsilon_1^2/2 \right) \right] \\ &\geq 1 - \delta \end{aligned} \tag{54}$$



## Lemma B.5

$$\begin{aligned} \mathbb{P} \left[ \overline{\mathcal{A}}^{h+1} \left( \epsilon_1^2/2 \right) \Rightarrow \overline{\mathcal{D}}^h \left( \epsilon_1 + \sqrt{\frac{2 \log \frac{6}{\delta_4}}{d_h}} \right) \right] &\geq 1 - \delta_4, \\ \forall \epsilon_1 \in [0, 1], \delta_4 \in (0, 1). \end{aligned} \tag{55}$$

# Proof of Theorem B.2 (17/47)

## Lemma B.8

Define  $\mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}') = \begin{bmatrix} \mathbf{g}^{(h)}(\mathbf{x})^\top \mathbf{g}^{(h)}(\mathbf{x}) & \mathbf{g}^{(h)}(\mathbf{x})^\top \mathbf{g}^{(h)}(\mathbf{x}') \\ \mathbf{g}^{(h)}(\mathbf{x}')^\top \mathbf{g}^{(h)}(\mathbf{x}) & \mathbf{g}^{(h)}(\mathbf{x}')^\top \mathbf{g}^{(h)}(\mathbf{x}') \end{bmatrix}$ , we have for every  $1 \leq h \leq L$ ,

$$\begin{aligned} \|\mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}') - \mathbf{\Lambda}^{(h)}(\mathbf{x}, \mathbf{x}')\|_\infty &\leq \frac{\epsilon^2}{2} \\ \Rightarrow \left| t_{\dot{\sigma}} \left( \mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}') \right) - t_{\dot{\sigma}} \left( \mathbf{\Lambda}^{(h)}(\mathbf{x}, \mathbf{x}') \right) \right| &\leq \epsilon, \forall 0 \leq \epsilon \leq 1. \end{aligned} \tag{56}$$



# Proof of Theorem B.2 (18/47)

## Proof of Lemma B.8.

For simplicity, we denote  $\mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}')$ ,  $\mathbf{\Lambda}^{(h)}(\mathbf{x}, \mathbf{x}')$  by  $\mathbf{G}, \mathbf{\Lambda}$  respectively. Since  $\dot{\sigma}(z) = \mathbf{1}[z \geq 0]$  is 0-homogeneous, we have

$$\begin{aligned} t_{\dot{\sigma}}(\mathbf{G}) &= \hat{t}_{\dot{\sigma}} \left( \frac{G_{12}}{\sqrt{G_{11}G_{22}}} \right) = \frac{1}{2} + \arcsin \frac{G_{12}}{\sqrt{G_{11}G_{22}}} \\ t_{\dot{\sigma}}(\mathbf{\Lambda}) &= \hat{t}_{\dot{\sigma}} \left( \frac{\Lambda_{12}}{\sqrt{\Lambda_{11}\Lambda_{22}}} \right) = \frac{1}{2} + \arcsin \frac{\Lambda_{12}}{\sqrt{\Lambda_{11}\Lambda_{22}}} = \frac{1}{2} + \arcsin \Lambda_{12} \end{aligned} \quad (57)$$

It is easy to verify that  $|\sqrt{G_{11}G_{22}} - 1| \leq \epsilon^2/2$ , and thus

$$\begin{aligned} \left| \frac{G_{12}}{\sqrt{G_{11}G_{22}}} - \Lambda_{12} \right| &\leq \left| \frac{G_{12}}{\sqrt{G_{11}G_{22}}} - \frac{\Lambda_{12}}{\sqrt{G_{11}G_{22}}} \right| + |\Lambda_{12}| \left| 1 - \frac{1}{\sqrt{G_{11}G_{22}}} \right| \\ &\leq \frac{\epsilon^2/2}{1 - \epsilon^2/2} + \frac{\epsilon^2/2}{1 - \epsilon^2/2} \leq 2\epsilon^2. \end{aligned} \quad (58)$$

# Proof of Theorem B.2 (19/47)

Proof of Lemma B.8.

Thus, by Lemma B.1

$$|t_{\dot{\sigma}}(\mathbf{G}) - t_{\dot{\sigma}}(\mathbf{\Lambda})| \leq \left| \frac{1}{2} + \arcsin \frac{G_{12}}{\sqrt{G_{11}G_{22}}} - \frac{1}{2} + \arcsin \Lambda_{12} \right| \leq \epsilon. \quad (59)$$



# Proof of Theorem B.2 (20/47)

## Lemma B.9

For any  $0 \leq h \leq L - 1$ , any fixed  $\{\mathbf{W}^{(i)}\}_{i=1}^h$ , w.p.  $1 - \delta$  over the randomness of  $\mathbf{W}^{(h+1)} \in \mathbb{R}^{d^{h+1} \times d^h}$ , we have

$$\left| 2 \frac{\text{tr}(\mathbf{D})}{d_h} - \hat{t}_{\dot{\sigma}} \left( \mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}') \right) \right| < \sqrt{\frac{2 \log \frac{2}{\delta}}{d_h}}. \quad (60)$$

## Proof of Lemma B.9.

Notice that  $\mathbb{E} 2 \frac{\text{tr}(\mathbf{D})}{d_h} = \hat{t}_{\dot{\sigma}} \left( \mathbf{G}^{(h)}(\mathbf{x}, \mathbf{x}') \right)$ , the proof is completed by Chernoff Bound. □

# Proof of Theorem B.2 (21/47)

## Proof of Lemma B.5.

Note that  $\dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') = t_{\sigma'} \left( \Sigma^{(h)} \Big|_{\mathbf{x}, \mathbf{x}'} \right) = \hat{t}_{\sigma'} \left( \Lambda^{(h)}(\mathbf{x}, \mathbf{x}') \right)$ .

Combining Lemma B.8 and Lemma B.9, we have for any  $(\mathbf{x}, \mathbf{x}')$ ,

$$\mathbb{P} \left[ \mathcal{D}^h \left( \mathbf{x}, \mathbf{x}', \epsilon_1 + \sqrt{\frac{2 \log \frac{6}{\delta}}{d_h}} \right) \mid \mathcal{A}^{h+1} \left( \mathbf{x}, \mathbf{x}', \epsilon_1^2/2 \right) \right] \geq 1 - \frac{\delta}{3}.$$

# Proof of Theorem B.2 (22/47)

## Proof of Lemma B.5.

Taking union bound over  $(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')$  for the choice of  $(\mathbf{x}, \mathbf{x}')$ , we have

$$\begin{aligned} & \mathbb{P} \left[ \overline{\mathcal{A}}^{h+1} \left( \epsilon_1^2/2 \right) \Rightarrow \overline{\mathcal{D}}^h \left( \epsilon_1 + \sqrt{\frac{2 \log \frac{6}{\delta}}{d_h}} \right) \right] \\ & \geq \mathbb{P} \left[ \overline{\mathcal{D}}^h \left( \epsilon_1 + \sqrt{\frac{2 \log \frac{6}{\delta}}{d_h}} \right) \mid \overline{\mathcal{A}}^{h+1} \left( \epsilon_1^2/2 \right) \right] \geq 1 - \delta \end{aligned} \tag{61}$$



# Proof of Theorem B.2 (23/47)

## Lemma B.7

There exists constant  $C, C' \in \mathbb{R}$ , for any  $\epsilon_2, \epsilon_3, \epsilon_4 \in [0, 1]$ , we have

$$\begin{aligned} & \mathbb{P} \left[ \overline{\mathcal{A}}^L (\epsilon_1^2/2) \wedge \overline{\mathcal{B}}^{h+1} (\epsilon_2) \wedge \overline{\mathcal{C}} (\epsilon_3) \wedge \overline{\mathcal{D}}^h (\epsilon_4) \right. \\ & \quad \left. \implies \overline{\mathcal{B}}^h \left( \epsilon_2 + \frac{C' \epsilon_3}{\sqrt{d_h}} + 2\epsilon_4 + C \sqrt{\frac{\log \frac{1}{\delta_2}}{d_h}} \right) \right] \\ & \geq 1 - \delta_2. \end{aligned} \tag{62}$$

# Proof of Theorem B.2 (24/47)

The proof of Lemma B.7 is based on the following 3 claims, Claim B.1, B.2 and B.3.

## Claim B.1

If  $\overline{\mathcal{A}}^L (\epsilon_1^2/2) \wedge \overline{\mathcal{B}}^{h+1} (\epsilon_2) \wedge \overline{\mathcal{C}} (\epsilon_3) \wedge \overline{\mathcal{D}}^h (\epsilon_4)$ , then we have

$$\left| \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \right\rangle - \prod_{h'=h}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| \leq \epsilon_2 + 2\epsilon_4. \quad (63)$$

# Proof of Theorem B.2 (25/47)

## Proof of Claim B.1.

$$\begin{aligned}
 & \left| \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \right\rangle - \prod_{h'=h}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| \\
 & \leq \left| \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} - \dot{\Sigma}^{(h)}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| \cdot \left| \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \right\rangle \right| \\
 & + \left| \dot{\Sigma}^{(h)}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| \cdot \left| \left\langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \right\rangle - \prod_{h'=h+1}^L \dot{\Sigma}^{(h')}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \right| \\
 & \leq 2\epsilon_4 + \epsilon_2
 \end{aligned} \tag{64}$$

□



## Proof of Theorem B.2 (26/47)

For any fixed  $h$ , let  $\mathbf{G} = [\mathbf{g}^{(h)}(\mathbf{x})\mathbf{g}^{(h)}(\mathbf{x}')]^T$ ,

### Claim B.2

*w.p.  $\geq 1 - \frac{\delta_2}{2}$ , if  $\overline{\mathcal{A}}^L(\epsilon_1^2/2) \wedge \overline{\mathcal{B}}^{h+1}(\epsilon_2) \wedge \overline{\mathcal{C}}(\epsilon_3) \wedge \overline{\mathcal{D}}^h(\epsilon_4)$ , then we have for any  $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,*

$$\left| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_{\mathbf{G}}^\perp \mathbf{D} \Pi_{\mathbf{G}}^\perp (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) - \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \rangle \right| \leq 16 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}}. \quad (65)$$

*As a by-product, for any  $\mathbf{x}^{(1)} \in \{\mathbf{x}, \mathbf{x}'\}$ , we have*

$$\sqrt{\frac{2}{d_h}} \|\mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_{\mathbf{G}}^\perp \mathbf{D}\| \leq 4 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}}. \quad (66)$$

# Proof of Theorem B.2 (27/47)

Lemma B.10 (Gaussian chaos of order 2,[BLM13])

Let  $\boldsymbol{\xi} \sim N(0, \mathbf{I}_n)$  be an  $n$ -dimensional unit gaussian random vector,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a symmetric matrix, then for any  $t > 0$ ,

$$\mathbb{P} \left[ \left| \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} - \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} \right| > 2 \|\mathbf{A}\|_F \sqrt{t} + 2 \|\mathbf{A}\|_2 t \right] \leq 2 \exp(-t).$$

Or,

$$\mathbb{P} \left[ \left| \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} - \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} \right| > t \right] \leq 2 \exp \left( - \frac{t^2}{4(\|\mathbf{A}\|_F^2) + \|\mathbf{A}\|_2 t} \right).$$

# Proof of Theorem B.2 (28/47)

## Proof of Claim B.2.

It suffices to prove this claim conditioned on every possible realization of

$$\{\mathbf{b}^{(h+1)}(\mathbf{x}^{(1)}), \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}), \mathbf{f}^{(h)}(\mathbf{x}^{(1)}), \mathbf{f}^{(h)}(\mathbf{x}^{(2)})\}.$$

Recall that  $\mathbf{G} = [\mathbf{g}^{(h)}(\mathbf{x}^{(1)})\mathbf{g}^{(h)}(\mathbf{x}^{(2)})]$ , we further define

$\mathbf{F} = [\mathbf{f}^{(h)}(\mathbf{x}^{(1)})\mathbf{f}^{(h)}(\mathbf{x}^{(2)})]$ . Applying Lemma B.3 on each row of  $\mathbf{W}^{h+1}$ , we have

$$\mathbf{W}^{(h+1)}\Pi_{\mathbf{G}}^{\perp} \stackrel{\text{d}}{=}_{\mathbf{F}=\mathbf{W}^{(h+1)}\mathbf{G}} \widetilde{\mathbf{W}}\Pi_{\mathbf{G}}^{\perp}, \quad (67)$$

where  $\widetilde{\mathbf{W}}$  is an iid copy of  $\mathbf{W}^{(h+1)}$ .

# Proof of Theorem B.2 (29/47)

## Proof of Claim B.2.

Note that  $[\mathbf{b}^{(h+1)}(\mathbf{x})^\top \widetilde{\mathbf{W}} \quad \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \widetilde{\mathbf{W}}]^\top \in \mathbb{R}^{2d_h}$  follows a joint zero-mean gaussian distribution with covariance matrix

$$\Sigma = \begin{bmatrix} b_{11}\mathbf{I}_n & b_{12}\mathbf{I}_n \\ b_{21}\mathbf{I}_n & b_{22}\mathbf{I}_n \end{bmatrix}, \text{ where } b_{ij} = \mathbf{b}^{(h)}(\mathbf{x}^{(i)})^\top \mathbf{b}^{(h)}(\mathbf{x}^{(j)}), \text{ for } i, j = 1, 2. \text{ In}$$

other words, there exists  $\mathbf{M} \in \mathbb{R}^{2d_h \times 2d_h}$ , s.t.  $\mathbf{M}\mathbf{M}^\top = \Sigma$ , and

$$[\mathbf{b}^{(h+1)}(\mathbf{x})^\top \widetilde{\mathbf{W}} \quad \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \widetilde{\mathbf{W}}]^\top \stackrel{\text{d}}{=} \mathbf{M}\boldsymbol{\xi},$$

where  $\boldsymbol{\xi} \sim N(\mathbf{0}, \mathbf{I}_{2d_h})$ .

# Proof of Theorem B.2 (30/47)

## Proof of Claim B.2.

Thus conditioned on  $\{\mathbf{b}^{(h+1)}(\mathbf{x}^{(1)}), \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}), \mathbf{g}^{(h)}(\mathbf{x}^{(1)}), \mathbf{g}^{(h)}(\mathbf{x}^{(2)})\}$ , we have

$$\begin{aligned}
 & \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp \mathbf{D} \Pi_G^\perp \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \\
 & \stackrel{\text{d}}{=} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \widetilde{\mathbf{W}} \Pi_G^\perp \mathbf{D} \Pi_G^\perp \left( \widetilde{\mathbf{W}} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \\
 & \stackrel{\text{d}}{=} \left( \begin{bmatrix} \mathbf{I}_{d_h} & \mathbf{0} \end{bmatrix} M \boldsymbol{\xi} \right)^\top \Pi_G^\perp \mathbf{D} \Pi_G^\perp \left( \begin{bmatrix} \mathbf{0} & \mathbf{I}_{d_h} \end{bmatrix} M \boldsymbol{\xi} \right) \\
 & \stackrel{\text{d}}{=} \frac{1}{2} \boldsymbol{\xi}^\top M^\top \begin{bmatrix} \mathbf{0} & \Pi_G^\perp \mathbf{D} \Pi_G^\perp \\ \Pi_G^\perp \mathbf{D} \Pi_G^\perp & \mathbf{0} \end{bmatrix} M \boldsymbol{\xi}.
 \end{aligned}$$

# Proof of Theorem B.2 (31/47)

## Proof of Claim B.2.

Now we are ready to prove Claim B.1 by applying Lemma B.10. Let

$$\mathbf{A} = \frac{1}{2} \mathbf{M}^\top \begin{bmatrix} \mathbf{0} & \Pi_G^\perp \mathbf{D} \Pi_G^\perp \\ \Pi_G^\perp \mathbf{D} \Pi_G^\perp & \mathbf{0} \end{bmatrix} \mathbf{M}, \text{ we have}$$

$$\begin{aligned} \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} &= \text{tr}(\mathbf{A}) = \frac{1}{2} \text{tr} \left( \begin{bmatrix} \mathbf{0} & \Pi_G^\perp \mathbf{D} \Pi_G^\perp \\ \Pi_G^\perp \mathbf{D} \Pi_G^\perp & \mathbf{0} \end{bmatrix} \boldsymbol{\Sigma} \right) \\ &= b_{12} \text{tr} \left( \Pi_G^\perp \mathbf{D} \Pi_G^\perp \mathbf{I}_n \right) = b_{12} \text{tr} \left( \mathbf{D} \Pi_G^\perp \right). \end{aligned} \quad (68)$$

Note that by definition  $\Pi_G^\perp = \mathbf{I}_{d_h} - \Pi_G$ , and  $\text{rank}(\Pi_G) \leq 2$ , we have

$$\begin{aligned} \text{tr} \left( \mathbf{D} \Pi_G^\perp \right) &= \text{tr} \left( \mathbf{D} (\mathbf{I} - \Pi_G) \right) = \text{tr}(\mathbf{D}) - \text{tr}(\mathbf{D} \Pi_G) \\ &= \text{tr}(\mathbf{D}) - \text{tr}(\Pi_G \mathbf{D} \Pi_G). \end{aligned} \quad (69)$$

# Proof of Theorem B.2 (32/47)

## Proof of Claim B.2.

Since  $\mathbf{0} \preceq \mathbf{D} \preceq \mathbf{I}_{d_h}$ , we have  $0 \leq \text{tr}(\Pi_G \mathbf{D} \Pi_G) \leq 2$ , and thus  $b_{12}(\text{tr}(\mathbf{D}) - 2) \leq \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} \leq b_{12} \text{tr}(\mathbf{D})$ . For the upper bound of spectrum, note that  $\|\mathbf{M}\|_2^2 = \|\boldsymbol{\Sigma}\|_2 = \left\| \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \right\|_2 \leq b_{11} + b_{12}$ , and  $\mathbf{0} \preceq \Pi_G^\perp, \mathbf{D} \preceq \mathbf{I}_{d_h}$ , we have

$$\begin{aligned} \|\mathbf{A}\|_2 &\leq \frac{1}{2} \|\mathbf{M}\|_2^2 \|\Pi_G^\perp \mathbf{D} \Pi_G^\perp\|_2 \leq \frac{1}{2} \|\mathbf{M}\|_2^2 \|\Pi_G^\perp\|_2 \|\mathbf{D}\|_2 \|\Pi_G^\perp\|_2 \\ &\leq \frac{b_{11} + b_{12}}{2} \leq \sqrt{2}, \end{aligned} \tag{70}$$

and

$$\|\mathbf{A}\|_F \leq \sqrt{2d_h} \|\mathbf{A}\|_2 = \frac{\sqrt{2d_h}(b_{11} + b_{22})}{2} \leq 2\sqrt{d_h}.$$

# Proof of Theorem B.2 (33/47)

## Proof of Claim B.2.

Thus by Lemma B.10 with  $t = \log \frac{6}{\delta_2}$  we have w.p.  $1 - \frac{\delta_2}{6}$ ,

$$\begin{aligned} \frac{1}{d_h} \left| \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} - \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} \right| &\leq \frac{1}{d_h} \left( 2 \|\mathbf{A}\|_F \sqrt{t} + 2 \|\mathbf{A}\|_2 t \right) \\ &= 4 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} + 2 \sqrt{2} \frac{\log \frac{6}{\delta_2}}{d_h} .. \end{aligned} \tag{71}$$



# Proof of Theorem B.2 (34/47)

## Proof of Claim B.2.

Thus we have

$$\begin{aligned}
 & \left| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp \mathbf{D} \Pi_G^\perp (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right. \\
 & \quad \left. - \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \rangle \right| \\
 & \leq \frac{2}{d_h} \left| \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} - \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} \right| + \left| 2 \mathbb{E} \boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} - \frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \rangle \right| \\
 & \leq 8 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} + 4 \sqrt{2} \frac{\log \frac{6}{\delta_2}}{d_h} + \frac{4b_{12}}{d_h}
 \end{aligned} \tag{72}$$

# Proof of Theorem B.2 (35/47)

Proof of Claim B.2.

$$\begin{aligned}
 &\leq 14 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} + \frac{4(1 + \epsilon_2)}{d_h} \quad (2\sqrt{2} \leq 3 \wedge \log \frac{6}{\delta_2} \leq d_h) \\
 &\leq 16 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} \quad (\epsilon_2 \leq 1 \wedge \sqrt{d_h \log 6} \geq 4).
 \end{aligned} \tag{73}$$

# Proof of Theorem B.2 (36/47)

## Proof of Claim B.2.

The main part of the claim is completed by taking union bound over  $(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')$ . For the by-product, let  $\mathbf{x}^{(2)} = \mathbf{x}^{(1)}$ , and we have

$$\begin{aligned}
 & \sqrt{\frac{2}{d_h}} \|\mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp \mathbf{D}\| \\
 & \leq \sqrt{\left| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp \mathbf{D} \Pi_G^\perp (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right|} \\
 & \leq \sqrt{\frac{2 \operatorname{tr}(\mathbf{D})}{d_h} \langle \mathbf{b}^{(h)}(\mathbf{x}^{(2)}), \mathbf{b}^{(h)}(\mathbf{x}^{(1)}) \rangle + \left( 16 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} \right)^2}
 \end{aligned} \tag{74}$$

# Proof of Theorem B.2 (37/47)

Proof of Claim B.2.

$$\begin{aligned} &\leq \sqrt{4 + \left(16 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}}\right)^2} \\ &\leq 2 + 4 \sqrt{\frac{\log \frac{6}{\delta_2}}{d_h}} \leq 6 \quad (\log \frac{6}{\delta_2} \leq d_h) \end{aligned} \tag{75}$$



# Proof of Theorem B.2 (38/47)

## Claim B.3

*w.p.  $\geq 1 - \frac{\delta_2}{2}$ , if  $\overline{\mathcal{A}}^L(\epsilon_1^2/2) \wedge \overline{\mathcal{B}}^{h+1}(\epsilon_2) \wedge \overline{\mathcal{C}}(\epsilon_3) \wedge \overline{\mathcal{D}}^h(\epsilon_4)$ , then*

$$\begin{aligned}\|\Pi_G \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x})\| &\leq 2\sqrt{\log \frac{8}{\delta_2}} + \sqrt{2}\epsilon_3, \\ \|\Pi_G \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}(\mathbf{x}')\| &\leq 2\sqrt{\log \frac{8}{\delta_2}} + \sqrt{2}\epsilon_3.\end{aligned}\tag{76}$$

# Proof of Theorem B.2 (39/47)

## Proof of Claim B.3.

It suffices to prove the claim for  $\mathbf{x}$ . We will denote  $\mathbf{x}$  by  $\mathbf{x}$ ,  $\mathbf{g}^{(h)}(\mathbf{x})$  by  $\mathbf{g}^{(h)}$  and  $\mathbf{b}^{(h+1)}(\mathbf{x})$  by  $\mathbf{b}^{(h+1)}$ . We also define  $\Pi_{\mathbf{g}}$  as  $\mathbf{g}\mathbf{g}^\top$ , and  $\Pi_{\mathbf{G}/\mathbf{g}} = \Pi_{\mathbf{G}} - \Pi_{\mathbf{g}}$ . Clearly,  $\Pi_{\mathbf{G}/\mathbf{g}}$  is still a projection matrix of rank 0 or 1.

Since  $\|\Pi_{\mathbf{G}} (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x})\| \leq \|\Pi_{\mathbf{g}} (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}\| + \|\Pi_{\mathbf{G}/\mathbf{g}} (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}\|$ , it suffices to bound these two terms separately.

# Proof of Theorem B.2 (40/47)

## Proof of Claim B.3.

Recall  $\mathbf{b}^{(h+1)}$  is defined as the gradient of  $f(\boldsymbol{\theta}, \mathbf{x})$  with respect to the pre-activation of layer  $h+1$ ,  $\mathbf{f}^{h+1}$ , thus if we view  $g$  as a function  $\mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)}$ , by the rule of back propagation, we have

$$\frac{\partial g(\mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)})}{\partial \mathbf{g}^{(h)}} = (\mathbf{b}^{(h+1)})^\top \mathbf{W}^{(h+1)}.$$

Note that  $\text{relu}$  is 1-homogeneous, namely  $\forall \lambda \in \mathbb{R}, \sigma(\lambda z) = \lambda \sigma(z)$ , the whole network is also 1-homogeneous in  $\mathbf{g}^{(h)}$ . In other words, we have

# Proof of Theorem B.2 (41/47)

## Proof of Claim B.3.

$$\begin{aligned}
 & f(\mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)}) \\
 &= \left. \frac{\partial f(\lambda \mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)})}{\partial \lambda} \right|_{\lambda=1} \\
 &= \left\langle \left. \frac{\partial f(\lambda \mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)})}{\partial \lambda \mathbf{g}^{(h)}} \right|_{\lambda=1}, \left. \frac{\partial \lambda \mathbf{g}^{(h)}}{\partial \lambda} \right|_{\lambda=1} \right\rangle \quad (77) \\
 &= \left\langle \frac{\partial g(\mathbf{g}^{(h)}, \mathbf{W}^{(h+1)}, \dots, \mathbf{W}^{(L+1)})}{\partial \mathbf{g}^{(h)}}, \mathbf{g}^{(h)} \right\rangle \\
 &= (\mathbf{g}^{(h)})^\top (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}
 \end{aligned}$$



# Proof of Theorem B.2 (42/47)

## Proof of Claim B.3.

By definition of  $\Pi_g$ , we have

$$\begin{aligned}
 \|\Pi_g \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}\| &= \left\| \frac{\mathbf{g}^{(h)} (\mathbf{g}^{(h)})^\top}{\|\mathbf{g}^{(h)}\|^2} \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)} \right\| \\
 &= \left| \frac{(\mathbf{g}^{(h)})^\top}{\|\mathbf{g}^{(h)}\|} \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)} \right| \\
 &= \frac{|f(\boldsymbol{\theta}, \mathbf{x})|}{\|\mathbf{g}^{(h)}\|}.
 \end{aligned} \tag{78}$$

# Proof of Theorem B.2 (43/47)

## Proof of Claim B.3.

Note that  $\mathbf{g}^{(h)}(x^{(0)})^\top \mathbf{g}^{(h)}(\mathbf{x}) \geq 1 - \epsilon_1^2/2 \geq \frac{1}{2}$ , we have

$$\|\Pi_{\mathbf{g}} \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}\| = \frac{|f(\boldsymbol{\theta}, \mathbf{x})|}{\|\mathbf{g}^{(h)}\|} \leq \sqrt{2}\epsilon_3.$$

For the second term  $\Pi_{\mathbf{G}/\mathbf{g}} \left( \mathbf{W}^{(h+1)} \right)^\top \mathbf{b}^{(h+1)}$ , note that conditioned on  $\mathbf{g}^{(h)}$ ,  $\mathbf{f}^h = \frac{1}{\sqrt{d_{h+1}}} \mathbf{W}^{(h+1)} \mathbf{g}^{(h)}$  and all  $\{\mathbf{W}^{(h)}\}_{h'}^{L+1}$  (thus  $\mathbf{b}^{(h+1)}$ ), by

Lemma B.3,  $\Pi_{\mathbf{G}/\mathbf{g}}(\mathbf{W}^{(h+1)}) \stackrel{\text{d}}{=} \Pi_{\mathbf{G}/\mathbf{g}} \widetilde{\mathbf{W}}$ , where  $\widetilde{\mathbf{W}}$  is an iid copy of  $\mathbf{W}^{(h+1)}$ .

# Proof of Theorem B.2 (44/47)

## Proof of Claim B.3.

Thus if  $\text{rank}(\mathbf{\Pi}_{G/g}) = 1$ , suppose  $\mathbf{\Pi}_{G/g} = \mathbf{u}\mathbf{u}^\top$  for some unit vector  $\mathbf{u}$ , we have

$$\begin{aligned}\|\mathbf{\Pi}_{G/g} \left(\mathbf{W}^{(h+1)}\right)^\top \mathbf{b}^{(h+1)}\| &= \left| \mathbf{u}^\top \left(\mathbf{W}^{(h+1)}\right)^\top \mathbf{b}^{(h+1)} \right| \\ &\stackrel{\text{d}}{=} \left| \mathbf{u}^\top \left(\widetilde{\mathbf{W}}\right)^\top \mathbf{b}^{(h+1)} \right| \\ &\stackrel{\text{d}}{=} |t|,\end{aligned}\tag{79}$$

where  $t \sim N(0, \|\mathbf{b}^{(h+1)}\|)$ . Hence w.p.  $\geq 1 - \delta_2/4$  over the randomness of  $\mathbf{W}^{(L)}$ ,  $\|\mathbf{\Pi}_{G/g} \left(\mathbf{W}^{(h+1)}\right)^\top \mathbf{b}^{(h+1)}\| \leq \sqrt{2 \log \frac{8}{\delta_2}} \|\mathbf{b}^{(h+1)}\| \leq \sqrt{2 \log \frac{8}{\delta_2}} \leq 2\sqrt{\log \frac{8}{\delta_2}} \ (\epsilon_2 < 1)$ .

# Proof of Theorem B.2 (45/47)

## Proof of Claim B.3.

If  $\text{rank}(\mathbf{\Pi}_{G/g}) = 0$ , then  $\|\mathbf{\Pi}_{G/g} (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}\| = 0 < 2\sqrt{\log \frac{8}{\delta_2}}$ .

Thus w.p.  $\geq 1 - \frac{\delta_2}{4}$ ,

$$\begin{aligned} & \|\mathbf{\Pi}_G (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x})\| \\ & \leq \|\mathbf{\Pi}_g (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}\| + \|\mathbf{\Pi}_{G/g} (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}\| \quad (80) \\ & \leq 2\sqrt{\log \frac{8}{\delta_2}} + \sqrt{2}\epsilon_3. \end{aligned}$$

Thus by assumption  $\log \frac{8}{\delta_2} \leq d_h$ , we have

$$2\sqrt{\log \frac{8}{\delta_2}} + \sqrt{2}\epsilon_3 \leq 2\sqrt{d_h} + \sqrt{2} \leq 3\sqrt{2d_h}.$$

□

# Proof of Theorem B.2 (46/47)

Wrapping things up, by combining Claim B.2 and Claim B.3, we have w.p.  $\geq 1 - \delta_2$ , for any pair of  $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,

$$\begin{aligned}
 & \left| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top (\mathbf{W}^{(h+1)})^\top D^{(h)}(\mathbf{x}^{(1)}) D^{(h)}(\mathbf{x}^{(2)}) (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) - \right. \\
 & \left. \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp D \Pi_G^\perp (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right| \\
 & \leq \left\| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G D \right\| \cdot \left\| D \Pi_G^\perp (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right\| \\
 & + \left\| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G^\perp D \right\| \cdot \left\| D \Pi_G (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right\| \\
 & + \left\| \frac{2}{d_h} \mathbf{b}^{(h+1)}(\mathbf{x}^{(1)})^\top \mathbf{W}^{(h+1)} \Pi_G \right\| \cdot \left\| \Pi_G (\mathbf{W}^{(h+1)})^\top \mathbf{b}^{(h+1)}(\mathbf{x}^{(2)}) \right\|
 \end{aligned} \tag{81}$$

# Proof of Theorem B.2 (47/47)

$$\begin{aligned}
 &\leq \left( 12\sqrt{2}\sqrt{\frac{\ln \frac{8}{\delta_2}}{d_h}} + 12\epsilon_3 \right) + \left( 12\sqrt{2}\sqrt{\frac{\ln \frac{8}{\delta_2}}{d_h}} + 12\epsilon_3 \right) + \left( 12\sqrt{2}\sqrt{\frac{\ln \frac{8}{\delta_2}}{d_h}} + 12\epsilon_3 \right) \\
 &= 36\sqrt{\frac{2 \ln \frac{8}{\delta_2}}{d_h}} + 36\epsilon_3.
 \end{aligned}
 \tag{82}$$

Using Equation (81) together with Claim B.1 and Claim B.2, we've finished the proof for Lemma B.7.

# Contents

## 1 Neural Tangent Kernels (NTK)

- Contruction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# Lemma 3.1 [ADH<sup>+</sup>19] (1/4)

## Lemma 3.1 (Evolution of outputs.)

- Recall in minimizing

$$\ell(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i)^2, \quad (83)$$

with infinitesimally small learning rate, one has

$$\frac{d\boldsymbol{\theta}(t)}{dt} = -\nabla \ell(\boldsymbol{\theta}(t)). \quad (84)$$



# Lemma 3.1 [ADH<sup>+</sup>19] (2/4)

## Lemma 3.1 (Evolution of outputs.)

- Let  $\mathbf{u}(t) = (f(\boldsymbol{\theta}(t), \mathbf{x}_i))_{i \in [n]} \in \mathbb{R}^n$  be the network outputs on all  $\mathbf{x}_i$ 's at time  $t$ , and  $\mathbf{y} = (y_i)_{i \in [n]}$  be the desired outputs.

Then  $\mathbf{u}(t)$  follows

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}(t) \cdot (\mathbf{u}(t) - \mathbf{y}), \quad (85)$$

where  $\mathbf{H}(t)$  is an  $n \times n$  positive semidefinite matrix whose  $(i, j)$ -th entry is  $\left\langle \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_j)}{\partial \boldsymbol{\theta}} \right\rangle$ :

## Lemma 3.1 [ADH<sup>+</sup>19] (3/4)

### Lemma 3.1 (Evolution of outputs.)

- For *over-parameterized* deep neural networks including (i) fully-connected neural networks with wide hidden layers and (ii) convolutional neural networks with many channels in hidden layers, with suitable parameterization and initialization schemes, it can be shown that the kernel  $\mathbf{H}(t)$  remains approximately unchanged during training, i.e.

$$\mathbf{H}(t) \approx \mathbf{H}(0). \quad (86)$$

Furthermore, as the hidden widths go to infinity, the (stochastic) kernel  $\mathbf{H}(0)$  at initialization converges to a deterministic limit  $\mathbf{H}^*$ . This gives rise to the *neural tangent kernel (NTK)*, as described in Equation (1), with  $[\mathbf{H}^*]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ . In Section ??, we describe the specific parameterizations for fully-connected and convolutional neural networks that give rise to the NTK phenomenon.

## Lemma 3.1 [ADH<sup>+</sup>19] (4/4)

### Lemma 3.1 (Evolution of outputs.)

Under approximations  $\mathbf{H}(t) \approx \mathbf{H}(0) \approx \mathbf{H}^*$ , if we work under the infinite width limit, Equation (90) can be replaced with the following *linear* dynamical system:

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}^* \cdot (\mathbf{u}(t) - \mathbf{y}).$$

This is the same as the dynamics of *kernel regression* under gradient flow, for which at time  $t \rightarrow \infty$  the final prediction function is (assuming  $\mathbf{u}(0) = \mathbf{0}$ )

$$f^*(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)) \cdot (\mathbf{H}^*)^{-1} \mathbf{y}. \quad (87)$$

Therefore, we can evaluate the performance of infinitely wide deep neural networks by looking at the kernel regression solution (87) using NTK. The problem boils down to the computation of NTK, which is the subject of Section ??.

# Proof of Lemma 3.1 [ADH<sup>+</sup>19] (1/2)

## Proof

The parameters  $\boldsymbol{\theta}$  evolve according to

$$\frac{d\boldsymbol{\theta}(t)}{dt} = -\nabla\ell(\boldsymbol{\theta}(t)) = -\sum_{i=1}^n (f(\boldsymbol{\theta}(t), \mathbf{x}_i) - y_i) \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{\partial \boldsymbol{\theta}}, \quad (88)$$

where  $t \geq 0$  is a continuous time index. Under Equation (88), the evolution of the network output  $f(\boldsymbol{\theta}(t), \mathbf{x}_i)$  can be written as

$$\begin{aligned} & \frac{d f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{d t} \\ &= -\sum_{j=1}^n (f(\boldsymbol{\theta}(t), \mathbf{x}_j) - y_j) \left\langle \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_j)}{\partial \boldsymbol{\theta}} \right\rangle, \quad \forall i \in [n]. \end{aligned} \quad (89)$$

# Proof of Lemma 3.1 [ADH<sup>+</sup>19] (2/2)

## Proof

Since  $\mathbf{u}(t) = (f(\boldsymbol{\theta}(t), \mathbf{x}_i))_{i \in [n]} \in \mathbb{R}^n$  is the network outputs on all  $\mathbf{x}_i$ 's at time  $t$ , and  $\mathbf{y} = (y_i)_{i \in [n]}$  is the desired outputs, Equation (89) can be written more compactly as

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}(t) \cdot (\mathbf{u}(t) - \mathbf{y}), \quad (90)$$

where  $\mathbf{H}(t) \in \mathbb{R}^{n \times n}$  is a kernel matrix defined as

$$\mathbf{H}(t)_{i,j} = \left\langle \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_j)}{\partial \boldsymbol{\theta}} \right\rangle, \quad (91)$$

$\forall i, j \in [n]$ .

# Equivalence Between NTK and NN (1/1)

Theorem 3.2 [ADH<sup>+</sup>19]

## Theorem 3.2 (Main theorem)

Suppose  $\sigma(z) = \max(0, z)$  ( $z \in \mathbb{R}$ ),  $1/\kappa = \text{poly}(1/\epsilon, \log(n/\delta))$  and  $d_1 = d_2 = \dots = d_L = m$  with  $m \geq \text{poly}(1/\kappa, L, 1/\lambda_0, n, \log(1/\delta))$ . Then for any  $\mathbf{x}_{te} \in \mathbb{R}^d$  with  $\|\mathbf{x}_{te}\| = 1$ , with probability at least  $1 - \delta$  over the random initialization, we have

$$|f_{nn}(\mathbf{x}_{te}) - f_{ntk}(\mathbf{x}_{te})| \leq \epsilon. \quad (92)$$

# Conclusions

- GD on DNN equivalent to Kernel GD wrt the NTK.
- NTK for MLPs ([JGH18, ADH<sup>+</sup>19]); NTK for convolutional nets (CNTK, [ADH<sup>+</sup>19]); NTK for graph nets (GNTK, [DHP<sup>+</sup>19])

# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]



# Notations and Setup

- Let  $G = (V, E)$  be a graph with node features  $\mathbf{h}_v \in \mathbb{R}^d$  for each  $v \in V$ .
- Denote the neighborhood of node  $v$  by  $\mathcal{N}(v)$ .
- Consider the graph classification task: given a set of graphs  $\{G_1, \dots, G_n\} \subseteq \mathcal{G}$  and their labels  $\{y_1, \dots, y_n\} \subseteq \mathcal{Y}$ , to learn to predict labels of unseen graphs.

# BLOCK Operation (1/2)

- A BLOCK operation aggregates features over a neighborhood  $\mathcal{N}(u) \cup \{u\}$  via, e.g., summation, and transforms the aggregated features with non-linearity, e.g., MLP or a fully-connected layer followed by ReLU.
- Denote the number of fully-connected layers in each BLOCK operation, i.e., the number of hidden layers of an MLP, by  $R$ .
- When  $R = 1$ , the BLOCK operation can be formulated as

$$\text{BLOCK}^{(\ell)}(u) = \sqrt{\frac{c_\sigma}{m}} \cdot \sigma \left( \mathbf{W}_\ell \cdot c_u \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v^{(\ell-1)} \right). \quad (93)$$

Here,  $\mathbf{W}_\ell$  are learnable weights, initialized as Gaussian random variables.  $\sigma$  is an activation function like ReLU.  $m$  is the output dimension of  $\mathbf{W}_\ell$ .

# BLOCK Operation (2/2)

- When the number of fully-connected layers  $R = 2$ , the BLOCK operation can be written as

$$\begin{aligned} & \text{BLOCK}^{(\ell)}(u) \\ &= \sqrt{\frac{c_\sigma}{m}} \sigma \left( \mathbf{W}_{\ell,2} \sqrt{\frac{c_\sigma}{m}} \cdot \sigma \left( \mathbf{W}_{\ell,1} \cdot c_u \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v^{(\ell-1)} \right) \right), \end{aligned} \quad (94)$$

where  $\mathbf{W}_{\ell,1}$  and  $\mathbf{W}_{\ell,2}$  are learnable weights. BLOCK operations can be defined similarly for  $R > 2$ .

# READOUT Operation

- To get the representation of an entire graph  $\mathbf{h}_G$  after  $L$  steps of aggregation, we take the summation over all node features, i.e.,

$$\mathbf{h}_G = \text{READOUT} \left( \left\{ \mathbf{h}_u^{(L)}, u \in V \right\} \right) = \sum_{u \in V} \mathbf{h}_u^{(L)}. \quad (95)$$

- There are more sophisticated READOUT operations than a simple summation. Jumping Knowledge Network (JK-Net) considers graph structures of different granularity, and aggregates graph features across all layers as

$$\mathbf{h}_G = \text{READOUT}^{\text{JK}} \left( \left\{ \mathbf{h}_u^{(\ell)}, u \in V, \ell \in [L] \right\} \right) = \sum_{u \in V} \left[ \mathbf{h}_u^{(0)}; \dots; \mathbf{h}_u^{(L)} \right]. \quad (96)$$

# Building GNNs using BLOCK and READOUT

- Most modern GNNs are constructed using the BLOCK operation and the READOUT operation. We denote the number of BLOCK operations (aggregation steps) in a GNN by  $L$ .
- For each  $\ell \in [L]$  and  $u \in V$ , we define  $\mathbf{h}_u^{(\ell)} = \text{BLOCK}^{(\ell)}(u)$ . The graph-level feature is then

$$\mathbf{h}_G = \text{READOUT} \left( \left\{ \mathbf{h}_u^{(L)}, u \in V \right\} \right) \quad (97)$$

or

$$\mathbf{h}_G = \text{READOUT}^{\text{JK}} \left( \left\{ \mathbf{h}_u^{(\ell)}, u \in V, \ell \in [L] \right\} \right), \quad (98)$$

depending on whether jumping knowledge is applied or not.

# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# GNTK: Problem Formulation

- Let  $f(\theta, G) \in \mathbb{R}$  be the output of the corresponding GNN under parameters  $\theta$  and input graph  $G$ , for two given graphs  $G$  and  $G'$ , to calculate the corresponding GNTK value, we need to calculate the expected value of

$$\left\langle \frac{\partial f(\theta, G)}{\partial \theta}, \frac{\partial f(\theta, G')}{\partial \theta} \right\rangle \quad (99)$$

in the limit that  $m \rightarrow \infty$  and  $\theta$  are all Gaussian random variables, which can be viewed as a Gaussian process.

- For each layer in the GNN, denote by  $\Sigma$  the covariance matrix of outputs of that layer and  $\dot{\Sigma}$  the covariance matrix corresponds to the derivative of that layer.
- These covariance matrices can be calculated via dynamic programming.



# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- **Formulae**
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# GNTK: Formulae

## Setup

- Given two graphs  $G = (V, E), G' = (V', E')$  with  $|V| = n, |V'| = n'$  and a GNN with  $L$  BLOCK operations and  $R$  fully-connected layers with ReLU activation in each BLOCK operation.
- First define the covariance matrix between input features of two input graphs  $G, G', \Sigma^{(0)}(G, G') \in \mathbb{R}^{n \times n'}$ .
- For two nodes  $u \in V$  and  $u' \in V', [\Sigma^{(0)}(G, G')]_{uu'}$  is defined to be  $\mathbf{h}_u^\top \mathbf{h}_{u'}$ , where  $\mathbf{h}_u$  and  $\mathbf{h}_{u'}$  are the input features of  $u \in V$  and  $u' \in V'$ .

# GNTK: Formulae for BLOCK Operation (1/3)

- A BLOCK operation in GNTK calculates a covariance matrix  $\Sigma_{(R)}^{(\ell)}(G, G') \in \mathbb{R}^{n \times n'}$  using  $\Sigma_{(R)}^{(\ell-1)}(G, G') \in \mathbb{R}^{n \times n'}$ , and calculates intermediate kernel values  $\Theta_{(r)}^{(\ell)}(G, G') \in \mathbb{R}^{n \times n'}$ .
- First perform a neighborhood aggregation operation

$$\begin{aligned}
 \left[ \Sigma_{(0)}^{(\ell)}(G, G') \right]_{uu'} &= c_u c_{u'} \sum_{v \in \mathcal{N}(u) \cup \{u\}} \sum_{v' \in \mathcal{N}(u') \cup \{u'\}} \left[ \Sigma_{(R)}^{(\ell-1)}(G, G') \right]_{vv'}, \\
 \left[ \Theta_{(0)}^{(\ell)}(G, G') \right]_{uu'} &= c_u c_{u'} \sum_{v \in \mathcal{N}(u) \cup \{u\}} \sum_{v' \in \mathcal{N}(u') \cup \{u'\}} \left[ \Theta_{(R)}^{(\ell-1)}(G, G') \right]_{vv'}.
 \end{aligned} \tag{100}$$

Here define  $\Sigma_{(R)}^{(0)}(G, G')$  and  $\Theta_{(R)}^{(0)}(G, G')$  as  $\Sigma^{(0)}(G, G')$ .

# GNTK: Formulae for BLOCK Operation (2/3)

- Next perform  $R$  transformations that correspond to the  $R$  fully-connected layers with ReLU activation.  $\sigma(z) = \max\{0, z\}$ , and  $\dot{\sigma}(z) = \mathbf{1}[z \geq 0]$ .
- For each  $r \in [R]$ , define
  - For  $u \in V, u' \in V'$ ,

$$\begin{aligned}
 & \left[ \mathbf{A}_{(r)}^{(\ell)}(G, G') \right]_{uu'} \\
 &= \begin{pmatrix} \left[ \boldsymbol{\Sigma}_{(r-1)}^{(\ell)}(G, G) \right]_{u,u} & \left[ \boldsymbol{\Sigma}_{(r-1)}^{(\ell)}(G, G') \right]_{uu'} \\ \left[ \boldsymbol{\Sigma}_{(r-1)}^{(\ell)}(G', G) \right]_{uu'} & \left[ \boldsymbol{\Sigma}_{(r-1)}^{(\ell)}(G', G') \right]_{u'u'} \end{pmatrix} \in \mathbb{R}^{2 \times 2}. \quad (101)
 \end{aligned}$$

# GNTK: Formulae for BLOCK Operation (3/3)

- For  $u \in V, u' \in V'$ ,

$$\begin{aligned} \left[ \Sigma_{(r)}^{(\ell)}(G, G') \right]_{uu'} &= c_{\sigma} \mathbb{E}_{(a,b) \sim \mathcal{N}(\mathbf{0}, [\mathbf{A}_{(r)}^{(\ell)}(G, G')]_{uu'})} [\sigma(a) \sigma(b)], \\ \left[ \dot{\Sigma}_{(r)}^{(\ell)}(G, G') \right]_{uu'} &= c_{\sigma} \mathbb{E}_{(a,b) \sim \mathcal{N}(\mathbf{0}, [\mathbf{A}_{(r)}^{(\ell)}(G, G')]_{uu'})} [\dot{\sigma}(a) \dot{\sigma}(b)]. \end{aligned} \quad (102)$$

- For  $u \in V, u' \in V'$ ,

$$\begin{aligned} \left[ \Theta_{(r)}^{(\ell)}(G, G') \right]_{uu'} &= \left[ \Theta_{(r-1)}^{(\ell)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(r)}^{(\ell)}(G, G') \right]_{uu'} \\ &\quad + \left[ \Sigma_{(r)}^{(\ell)}(G, G') \right]_{uu'}. \end{aligned} \quad (103)$$

# GNTK: Formulae for READOUT Operation

- Given these intermediate outputs, the final output of GNTK is

$$\Theta(G, G') = \begin{cases} \sum_{u \in V, u' \in V'} \left[ \Theta_{(R)}^{(\ell)}(G, G') \right]_{uu'} \text{ w/o jumping,} \\ \sum_{u \in V, u' \in V'} \left[ \sum_{\ell=0}^L \Theta_{(R)}^{(\ell)}(G, G') \right]_{uu'} \text{ w/ jumping.} \end{cases} \quad (104)$$

# GNTK: An Example

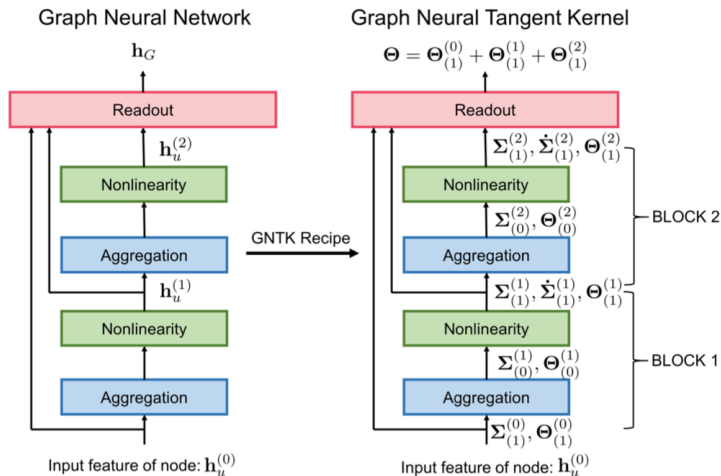


Figure: Exapmle of GNTK from a GNN.

# GNTK: An Example

- Consider a GNN with  $L = 2$  BLOCKS,  $R = 1$  fully-connected layer in each BLOCK, and jumping knowledge.
- For two graphs  $G$  and  $G'$ , first calculate

$$\left[\Theta_{(1)}^{(0)}(G, G')\right]_{uu'} = \left[\Sigma_{(1)}^{(0)}(G, G')\right]_{uu'} = \left[\Sigma^{(0)}(G, G')\right]_{uu'} = \mathbf{h}_u^\top \mathbf{h}_{u'}. \quad (105)$$

- Then follow the kernel formulae to calculate  $\Sigma_{(0)}^{(\ell)}, \Theta_{(0)}^{(\ell)}$  using  $\Sigma_{(R)}^{(\ell-1)}, \Theta_{(R)}^{(\ell-1)}$  (Aggregation) and calculate  $\Sigma_{(r)}^{(\ell)}, \dot{\Sigma}_{(r)}^{(\ell)}, \Theta_{(r)}^{(\ell)}$  using  $\Sigma_{(r-1)}^{(\ell)}, \Theta_{(r-1)}^{(\ell)}$  (Nonlinearity).
- The final output is

$$\Theta(G, G') = \sum_{u \in V, u' \in V'} \left[ \sum_{\ell=0}^L \Theta_{(R)}^{(\ell)}(G, G') \right]_{uu'}. \quad (106)$$



# Contents

## 1 Neural Tangent Kernels (NTK)

- Construction of the NTK due to [ADH<sup>+</sup>19]
- Convergence to the NTK at Initialization [ADH<sup>+</sup>19]
- Equivalence between NN and NTK

## 2 Graph Neural Nets (GNN)

## 3 Graph Neural Tangent Kernel (GNTK)

- Overview
- Formulae
- Theoretical Analyses
  - Proof of Theorem 4.2 [DHP<sup>+</sup>19]
  - Proof of Theorem 4.3 [DHP<sup>+</sup>19]

# GNTK: Generalization (1/2)

- Given  $n$  training data  $\{(G_i, y_i)\}_{i=1}^n$  drawn i.i.d. from the underlying distribution  $\mathcal{D}$ , where  $G_i$  is the  $i$ -th input graph and  $y_i$  is its label.
- Consider a GNN with a single BLOCK operation, followed by the READOUT operation (without jumping knowledge).
- Set  $c_u = \left( \left\| \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v \right\|_2 \right)^{-1}$ , and  $\Theta \in \mathbb{R}^{n \times n}$  to be the kernel matrix, where  $\Theta_{ij} = \Theta(G_i, G_j)$ .

# GNTK: Generalization (2/2)

- Assume the kernel matrix  $\Theta \in \mathbb{R}^{n \times n}$  is invertible.
- For a testing point  $G_{te}$ , the prediction of kernel regression using GNTK on this testing point is

$$f_{ker}(G_{te}) = [\Theta(G_{te}, G_1), \Theta(G_{te}, G_1), \dots, \Theta(G_{te}, G_n)]^\top \Theta^{-1} \mathbf{y}. \quad (107)$$

# GNTK Theorem 4.1 [DHP<sup>+</sup>19] (1/1)

## Theorem 4.1 (Generalization bound of kernel regression [BM02].)

Given  $n$  training data  $\{(G_i, y_i)\}_{i=1}^n$  drawn i.i.d. from the underlying distribution  $\mathcal{D}$ . Consider any loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$  that is 1-Lipschitz in the first argument such that  $\ell(y, y) = 0$ . With probability at least  $1 - \delta$ , the population loss of the GNTK predictor can be upper bounded by

$$\begin{aligned} L_{\mathcal{D}}(f_{\text{ker}}) &= \mathbf{E}_{(G, y) \sim \mathcal{D}} [\ell(f_{\text{ker}}(G), y)] \\ &= \leq O \left( \frac{\sqrt{\mathbf{y}^{\top} \Theta^{-1} \mathbf{y} \cdot \text{tr}(\Theta)}}{n} + \sqrt{\frac{\log(1/\delta)}{n}} \right). \end{aligned} \quad (108)$$

- Using the above theorem, if one can bound  $\mathbf{y}^{\top} \Theta^{-1} \mathbf{y}$  and  $\text{tr}(\Theta)$ , then one can obtain a sample complexity bound.

# GNTK Theorem 4.2 [DHP<sup>+</sup>19] (1/1)

Theorem 4.2 (An upper bound of  $\mathbf{y}^\top \Theta^{-1} \mathbf{y}$ )

For each  $i \in [n]$ , if the labels  $\{y_i\}_{i=1}^n$  satisfy

$$y_i = \alpha_1 \sum_{u \in V} (\bar{\mathbf{h}}_u^\top \beta_1) + \sum_{l=1}^{\infty} \alpha_{2l} \sum_{u \in V} (\bar{\mathbf{h}}_u^\top \beta_{2l})^{2l}, \quad (109)$$

where  $\bar{\mathbf{h}}_u = c_u \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v$ ,  $\alpha_1, \alpha_2, \alpha_4, \dots \in \mathbb{R}$ ,  $\beta_1, \beta_2, \beta_4, \dots \in \mathbb{R}^d$ , and  $G_i = (V, E)$ , then we have

$$\sqrt{\mathbf{y}^\top \Theta^{-1} \mathbf{y}} \leq 2|\alpha_1| \cdot \|\beta_1\|_2 + \sum_{l=1}^{\infty} \sqrt{2\pi}(2l-1)|\alpha_{2l}| \cdot \|\beta_{2l}\|_2^{2l}. \quad (110)$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (1/10)

For two graph  $G$  and  $G'$ , the GNTK kernel function that corresponds to the simple GNN can be described as

$$\begin{aligned} & \Theta(G, G') \\ &= \sum_{u \in V, u' \in V'} \left( \underbrace{\left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} + \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'}}_{\text{Theta\_1}} \right). \end{aligned} \quad (111)$$

Here, we have

$$\left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} = c_u c_{u'} \left( \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v \right)^\top \left( \sum_{v' \in \mathcal{N}(u') \cup \{u'\}} \mathbf{h}_{v'} \right) = \bar{\mathbf{h}}_u^\top \bar{\mathbf{h}}_{u'}. \quad (112)$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (2/10)

Recall that

$$\begin{aligned} \left[ \Sigma_{(r)}^{(\ell)}(G, G') \right]_{uu'} &= c_{\sigma} \mathbb{E}_{(a,b) \sim \mathcal{N}(\mathbf{0}, [A_{(r)}^{(\ell)}(G, G')]_{uu'})} [\sigma(a) \sigma(b)], \\ \left[ \dot{\Sigma}_{(r)}^{(\ell)}(G, G') \right]_{uu'} &= c_{\sigma} \mathbb{E}_{(a,b) \sim \mathcal{N}(\mathbf{0}, [A_{(r)}^{(\ell)}(G, G')]_{uu'})} [\dot{\sigma}(a) \dot{\sigma}(b)] \end{aligned} \quad (113)$$

and

$$\left[ A_{(r)}^{(\ell)}(G, G') \right]_{uu'} = \begin{pmatrix} \left[ \Sigma_{(r-1)}^{(\ell)}(G, G) \right]_{u,u} & \left[ \Sigma_{(r-1)}^{(\ell)}(G, G') \right]_{uu'} \\ \left[ \Sigma_{(r-1)}^{(\ell)}(G', G) \right]_{uu'} & \left[ \Sigma_{(r-1)}^{(\ell)}(G', G') \right]_{u'u'} \end{pmatrix} \in \mathbb{R}^{2 \times 2}. \quad (114)$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (3/10)

Since  $\sigma(z) = \max\{0, z\}$  is the ReLU activation function, and  $\dot{\sigma}(z) = \mathbb{1}[z \geq 0]$  is the derivative of the ReLU activation function, and  $\|\bar{\mathbf{h}}_u\|_2 = 1$  for all nodes  $u$ , by calculation, we have

$$\begin{aligned} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} &= \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right)}{2\pi}, \\ \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'} &= \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right) + \sqrt{1 - \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'}^2}}{2\pi}. \end{aligned} \quad (115)$$



# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (4/10)

Since

$$\arcsin(x) = \sum_{l=0}^{\infty} \frac{(2l-1)!!}{(2l)!!} \cdot \frac{x^{2l+1}}{2l+1}, \quad (116)$$

we have

$$\begin{aligned} & \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} \\ &= \frac{1}{4} \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} + \frac{1}{2\pi} \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \arcsin \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right) \\ &= \frac{1}{4} \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} + \frac{1}{2\pi} \sum_{l=1}^{\infty} \frac{(2l-3)!!}{(2l-2)!! \cdot (2l-1)} \cdot \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'}^{2l} \\ &= \frac{1}{4} \bar{\mathbf{h}}_u^\top \bar{\mathbf{h}}_{u'} + \frac{1}{2\pi} \sum_{l=1}^{\infty} \frac{(2l-3)!!}{(2l-2)!! \cdot (2l-1)} \cdot \left( \bar{\mathbf{h}}_u^\top \bar{\mathbf{h}}_{u'} \right)^{2l}. \end{aligned} \quad (117)$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (5/10)

Let  $\Phi^{(2l)}(\cdot)$  be the feature map of the polynomial kernel of degree  $2l$ , i.e.,

$$k^{(2l)}(\mathbf{x}, \mathbf{y}) = \left(\mathbf{x}^\top \mathbf{y}\right)^{2l} = \Phi^{(2l)}(\mathbf{x})^\top \Phi^{(2l)}(\mathbf{y}). \quad (118)$$

We have

$$\begin{aligned} & \left[\Sigma_{(0)}^{(1)}(G, G')\right]_{uu'} \left[\dot{\Sigma}_{(1)}^{(1)}(G, G')\right]_{uu'} \\ &= \frac{1}{4} \bar{\mathbf{h}}_u^\top \bar{\mathbf{h}}_{u'} + \frac{1}{2\pi} \sum_{l=1}^{\infty} \frac{(2l-3)!!}{(2l-2)!! \cdot (2l-1)} \cdot \left(\Phi^{(2l)}(\bar{\mathbf{h}}_u)\right)^\top \Phi^{(2l)}(\bar{\mathbf{h}}_{u'}). \end{aligned} \quad (119)$$

Let

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (6/10)

$$\Theta_1(G, G') = \sum_{u \in V, u' \in V'} \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'}, \quad (120)$$

we have

$$\begin{aligned} & \Theta_1(G, G') \\ &= \frac{1}{4} \left( \sum_{u \in V} \bar{\mathbf{h}}_u \right)^\top \left( \sum_{u' \in V'} \bar{\mathbf{h}}_{u'} \right) + \frac{1}{2\pi} \sum_{l=1}^{\infty} \frac{(2l-3)!!}{(2l-2)!! \cdot (2l-1)} \cdot \left( \sum_{u \in V} \Phi^{(2l)}(\bar{\mathbf{h}}_u) \right) \end{aligned} \quad (121)$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (7/10)

Since  $\Theta = \Theta_1 + \Theta_2$  where  $\Theta_2$  is a kernel matrix (and thus positive semi-definite), for any  $y \in \mathbb{R}^n$ , we have

$$y^\top \Theta^{-1} y \leq y^\top \Theta_1^{-1} y. \quad (122)$$

JY: I was not able to prove this yet.

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (8/10)

Recall that

$$y_i = \alpha_1 \sum_{u \in V} \left( \bar{\mathbf{h}}_u^\top \boldsymbol{\beta}_1 \right) + \sum_{l=1}^{\infty} \alpha_{2l} \sum_{u \in V} \left( \bar{\mathbf{h}}_u^\top \boldsymbol{\beta}_{2l} \right)^{2l}.$$

We rewrite

$$y_i = y_i^{(0)} + \sum_{l=1}^{\infty} y_i^{(2l)}, \quad (123)$$

where

$$y_i^{(0)} = \alpha_1 \left( \sum_{u \in V} \bar{\mathbf{h}}_u \right)^\top \boldsymbol{\beta}_1, \quad (124)$$

and for each  $l \geq 1$ ,

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (9/10)

$$\begin{aligned}
 y_i^{(2l)} &= \alpha_{2l} \sum_{u \in V} \left( \bar{\mathbf{h}}_u^\top \beta_{2l} \right)^{2l} = \alpha_{2l} \sum_{u \in V} \left( \Phi^{2l} \left( \bar{\mathbf{h}}_u \right) \right)^\top \Phi^{2l} (\beta_{2l}) \\
 &= \alpha_{2l} \left( \sum_{u \in V} \Phi^{2l} \left( \bar{\mathbf{h}}_u \right) \right)^\top \Phi^{2l} (\beta_{2l}).
 \end{aligned} \tag{125}$$

We have

$$\mathbf{y} = \mathbf{y}^{(0)} + \sum_{l=1}^{\infty} \mathbf{y}^{(2l)}. \tag{126}$$

Thus,

$$\sqrt{\mathbf{y}^\top \Theta^{-1} \mathbf{y}} \leq \sqrt{\mathbf{y}^\top \Theta_1^{-1} \mathbf{y}} \leq \sqrt{(\mathbf{y}^{(0)})^\top \Theta_1^{-1} \mathbf{y}^{(0)}} + \sum_{l=1}^{\infty} \sqrt{(\mathbf{y}^{(2l)})^\top \Theta_1^{-1} \mathbf{y}^{(2l)}}. \tag{127}$$

# Proof of Theorem 4.2 [DHP<sup>+</sup>19] (10/10)

When  $l = 0$ , we have

$$\sqrt{(\mathbf{y}^0)^\top \boldsymbol{\Theta}_1^{-1} \mathbf{y}^0} \leq 2|\alpha_1| \|\boldsymbol{\beta}_1\|_2. \quad (128)$$

When  $l \geq 1$ , we have

$$\sqrt{(\mathbf{y}^{2l})^\top \boldsymbol{\Theta}_1^{-1} \mathbf{y}^{2l}} \leq \sqrt{2\pi}(2l-1)|\alpha_{2l}| \left\| \boldsymbol{\Phi}^{2l}(\boldsymbol{\beta}_{2l}) \right\|_2. \quad (129)$$

Notice that

$$\left\| \boldsymbol{\Phi}^{2l}(\boldsymbol{\beta}_{2l}) \right\|_2^2 = \left( \boldsymbol{\Phi}^{2l}(\boldsymbol{\beta}_{2l}) \right)^\top \boldsymbol{\Phi}^{2l}(\boldsymbol{\beta}_{2l}) = \|\boldsymbol{\beta}_{2l}\|_2^{4l}. \quad (130)$$

Thus,

$$\sqrt{\mathbf{y}^\top \boldsymbol{\Theta}^{-1} \mathbf{y}} \leq 2|\alpha_1| \|\boldsymbol{\beta}_1\|_2 + \sum_{l=1}^{\infty} \sqrt{2\pi}(2l-1)|\alpha_{2l}| \|\boldsymbol{\beta}_{2l}\|_2^{2l}. \quad (131)$$

# GNTK Theorem 4.3 [DHP<sup>+</sup>19]

Theorem 4.3 (An upper bound of  $\text{tr}(\Theta)$ )

*If for all graphs  $G_i = (V_i, E_i)$  in the training set,  $|V_i|$  is upper bounded by  $\bar{V}$ , then  $\text{tr}(\Theta) \leq O(n\bar{V}^2)$ . Here  $n$  is the number of training samples.*

=



# Proof of Theorem 4.3 [DHP<sup>+</sup>19] (1/4)

Recall that

$$\begin{aligned} & \Theta(G, G') \\ &= \sum_{u \in V, u' \in V'} \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} + \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'} \right), \end{aligned} \quad (132)$$

where

$$\left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} = c_u c_{u'} \left( \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v \right)^\top \left( \sum_{v' \in \mathcal{N}(u') \cup \{u'\}} \mathbf{h}_{v'} \right) = \bar{\mathbf{h}}_u^\top \bar{\mathbf{h}}_{u'} \quad (133)$$

# Proof of Theorem 4.3 [DHP<sup>+</sup>19] (2/4)

and

$$\begin{aligned}
 \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} &= \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right)}{2\pi}, \\
 \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'} & \\
 &= \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right) + \sqrt{1 - \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'}^2}}{2\pi}.
 \end{aligned} \tag{134}$$

# Proof of Theorem 4.3 [DHP<sup>+</sup>19] (3/4)

Since for each node  $u$ ,  $\bar{\mathbf{h}}_u = c_u \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v$ , and  $c_u = \left( \left\| \sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v \right\|_2 \right)^{-1}$ , we have  $\|\bar{\mathbf{h}}_u\|_2 = 1$ . Moreover,

$$\left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} = \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right)}{2\pi} \leq 1/2 \quad (135)$$

and

$$\begin{aligned} & \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'} \\ &= \frac{\pi - \arccos \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \right) + \sqrt{1 - \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'}^2}}{2\pi} \leq \frac{1 + \pi}{2\pi} \leq 1, \end{aligned} \quad (136)$$

# Proof of Theorem 4.3 [DHP<sup>+</sup>19] (4/4)

we have

$$\begin{aligned}
 \Theta(G, G') &= \sum_{u \in V, u' \in V'} \left( \left[ \Sigma_{(0)}^{(1)}(G, G') \right]_{uu'} \left[ \dot{\Sigma}_{(1)}^{(1)}(G, G') \right]_{uu'} + \left[ \Sigma_{(1)}^{(1)}(G, G') \right]_{uu'} \right) \\
 &\leq \left( \frac{1}{2} + 1 \right) |V| |V'| \\
 &\leq 2|V| |V'|.
 \end{aligned} \tag{137}$$

Thus,

$$\text{tr}(\Theta) \leq 2n\bar{V}^2. \tag{138}$$

# GNTK: Remarks

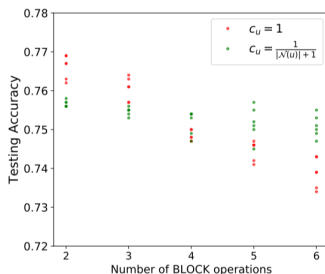
- Combining these three results, we know if

$$2|\alpha_1| \cdot \|\beta_1\|_2 + \sum_{l=1}^{\infty} \sqrt{2\pi}(2l-1)|\alpha_{2l}| \cdot \|\beta_{2l}\|_2^{2l} \quad (139)$$

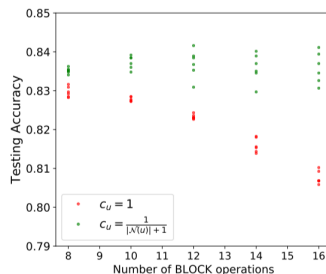
is bounded, and  $|V_i|$  is bounded for all graphs  $G_i = (V_i, E_i)$  in the training set, then the GNTK that corresponds to the simple GNN described above can learn functions of forms in (109), with polynomial number of samples.

# GNTK: Experiments I

- Test accuracy is correlated with the dataset and architecture.



(a) IMDBBINARY



(b) NCI1

Figure: Effect of number of BLOCK operations.

# GNTK: Experiments II

- Jump knowledge is expected to improve performance.
- The authors of [DHP<sup>+</sup>19] observed a 0.8% improvement in accuracy.

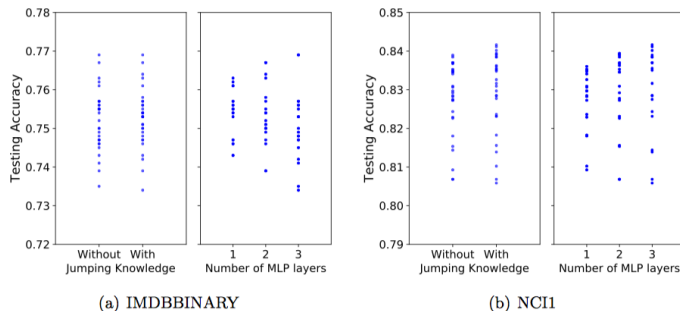


Figure: Effect of jump knowledge.

# References (1/2)



Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang, *On exact computation with an infinitely wide neural net*, arXiv preprint arXiv:1904.11955 (2019).



Stéphane Boucheron, Gábor Lugosi, and Pascal Massart, *Concentration inequalities: A nonasymptotic theory of independence*, Oxford university press, 2013.



Peter L Bartlett and Shahar Mendelson, *Rademacher and gaussian complexities: Risk bounds and structural results*, Journal of Machine Learning Research **3** (2002), no. Nov, 463–482.



Amit Daniely, Roy Frostig, and Yoram Singer, *Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity*, Advances In Neural Information Processing Systems, 2016, pp. 2253–2261.



Simon S. Du, Kangcheng Hou, Barnabás Póczos, Ruslan Salakhutdinov, Ruosong Wang, and Keyulu Xu, *Graph neural tangent kernel: Fusing graph neural networks with graph kernels*, CoRR **abs/1905.13192** (2019).



# References (2/2)



Arthur Jacot, Franck Gabriel, and Clément Hongler, *Neural tangent kernel: Convergence and generalization in neural networks*, Advances in neural information processing systems, 2018, pp. 8571–8580.



Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein, *Deep neural networks as gaussian processes*, arXiv preprint arXiv:1711.00165 (2017).



Greg Yang, *Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation*, arXiv preprint arXiv:1902.04760 (2019).