



# On Exact Computation with an Infinitely Wide Neural Net

Aurelien Lucchi

# 1. Recall NTK

# NTK

- Notation: Denote by  $f(\boldsymbol{\theta}, \mathbf{x}) \in \mathbb{R}$  the output of a neural network
  - $\boldsymbol{\theta} \in \mathbb{R}^N$  is all the parameters in the network
  - $\mathbf{x} \in \mathbb{R}^d$  is the input

# NTK

- Notation: Denote by  $f(\boldsymbol{\theta}, \mathbf{x}) \in \mathbb{R}$  the output of a neural network
  - $\boldsymbol{\theta} \in \mathbb{R}^N$  is all the parameters in the network
  - $\mathbf{x} \in \mathbb{R}^d$  is the input
- NTK is defined using the *gradient* of the output of the randomly initialized net with respect to its parameters, i.e.,

$$\text{ker}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{\theta} \sim \mathcal{W}} \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle. \quad (1)$$

- Generalization to convolutional neural nets: *Convolutional Neural Tangent Kernel (CNTK)*.

# Contributions

- **Theory:** Non-asymptotic proof that the NTK captures the behavior of a **fully-trained wide** neural net
  - weaker condition than previous proofs!
- **Empirically:**
  - New dynamic programming algorithm to compute CNTKs for ReLU activation (namely, to compute  $\ker(\mathbf{x}, \mathbf{x}')$ )
  - Evaluate the performance of fully-trained infinitely wide nets: performance on CIFAR-10 is within 5% of the performance in the finite case

## Gradient flow & NTK

- Training by minimizing the squared loss:

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

- Optimize using gradient flow:  $\frac{d\boldsymbol{\theta}(t)}{dt} = -\nabla \ell(\boldsymbol{\theta}(t))$

## Gradient flow & NTK

- Training by minimizing the squared loss:

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

- Optimize using gradient flow:  $\frac{d\boldsymbol{\theta}(t)}{dt} = -\nabla \ell(\boldsymbol{\theta}(t))$

### Lemma

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

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}(t) \cdot (\mathbf{u}(t) - \mathbf{y}). \quad (2)$$

## Limit behavior

- In the infinite limit width, matrix  $\mathbf{H}(t)$  remains *constant* during training i.e.  $\mathbf{H}(t) = \mathbf{H}(0)$ .
- Under a random initialization of  $\theta$ , as the width goes to infinity,  $\mathbf{H}(0)$  converges in probability to a deterministic kernel matrix  $\mathbf{H}^* = \text{ker}(\cdot, \cdot)$  (i.e. the NTK)
- $\implies$  Equation (2) becomes

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}^* \cdot (\mathbf{u}(t) - \mathbf{y}). \quad (3)$$



## Limit behavior

- In the infinite limit width, matrix  $\mathbf{H}(t)$  remains *constant* during training i.e.  $\mathbf{H}(t) = \mathbf{H}(0)$ .
- Under a random initialization of  $\theta$ , as the width goes to infinity,  $\mathbf{H}(0)$  converges in probability to a deterministic kernel matrix  $\mathbf{H}^* = \text{ker}(\cdot, \cdot)$  (i.e. the NTK)
- $\implies$  Equation (2) becomes

$$\frac{d\mathbf{u}(t)}{dt} = -\mathbf{H}^* \cdot (\mathbf{u}(t) - \mathbf{y}). \quad (3)$$

- Remark: Above dynamics is identical to the *dynamics of kernel regression under gradient flow*, for which at time  $t \rightarrow \infty$  the final prediction function is:

$$f^*(\mathbf{x}) = (\text{ker}(\mathbf{x}, \mathbf{x}_1), \dots, \text{ker}(\mathbf{x}, \mathbf{x}_n)) \cdot (\mathbf{H}^*)^{-1} \mathbf{y}. \quad (4)$$

## 2. Main results

# Network

- Define  $L$ -hidden-layer fully-connected neural network:

$$\begin{aligned} \mathbf{f}^{(h)}(\mathbf{x}) &= \mathbf{W}^{(h)} \mathbf{g}^{(h-1)}(\mathbf{x}) \in \mathbb{R}^{d_h} \\ \mathbf{g}^{(h)}(\mathbf{x}) &= \sqrt{\frac{c_\sigma}{d_h}} \sigma \left( \mathbf{f}^{(h)}(\mathbf{x}) \right) \in \mathbb{R}^{d_h}, \quad h = 1, 2, \dots, L, \end{aligned} \quad (5)$$

where  $\mathbf{W}^{(h)} \in \mathbb{R}^{d_h \times d_{h-1}}$ ,  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ ,  
 $c_\sigma = \left( \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[ \sigma(z)^2 \right] \right)^{-1}$ .

# Network

- Last layer of the neural network is

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

where  $\boldsymbol{\theta} = \left( \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L+1)} \right)$

# NTK (Lee et al. (2018))

- $\mathbf{f}^{(h)}(\mathbf{x}) = \mathbf{W}^{(h)} \mathbf{g}^{(h-1)}(\mathbf{x}) \in \mathbb{R}^{d_h}$
- Recall from Lee et al. (2018) that in the infinite width limit, the pre-activations  $\mathbf{f}^{(h)}(\mathbf{x})$  at every hidden layer  $h \in [L]$  has all its **coordinates tending to i.i.d. centered Gaussian processes** of covariance  $\Sigma^{(h-1)} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  defined recursively as:

$$\begin{aligned}\Sigma^{(0)}(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^\top \mathbf{x}', \\ \Lambda^{(h)}(\mathbf{x}, \mathbf{x}') &= \begin{pmatrix} \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \\ \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}) & \Sigma^{(h-1)}(\mathbf{x}', \mathbf{x}') \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \quad (6) \\ \Sigma^{(h)}(\mathbf{x}, \mathbf{x}') &= c_\sigma \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \Lambda^{(h)})} [\sigma(u) \sigma(v)].\end{aligned}$$

# NTK

- To give the formula of NTK, we also need to define a derivative covariance:

$$\dot{\Sigma}^{(h)}(\mathbf{x}, \mathbf{x}') = c_{\sigma} \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^{(h)})} [\dot{\sigma}(u) \dot{\sigma}(v)] . \quad (7)$$

- Final NTK expression for the fully-connected neural network is

$$\Theta^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{h=1}^{L+1} \left( \Sigma^{(h-1)}(\mathbf{x}, \mathbf{x}') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(\mathbf{x}, \mathbf{x}') \right), \quad (8)$$

where we let  $\dot{\Sigma}^{(L+1)}(\mathbf{x}, \mathbf{x}') = 1$  for convenience.

## NTK at initialization

Theorem (1) Convergence to the NTK at initialization)

Fix  $\epsilon > 0$  and  $\delta \in (0, 1)$ . Suppose  $\sigma(z) = \max(0, z)$  and  $\min_{h \in [L]} d_h \geq \Omega(\frac{L^6}{\epsilon^4} \log(L/\delta))$ . 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:

$$\left| \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle - \Theta^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L + 1)\epsilon.$$

## NTK at initialization

Theorem (1) Convergence to the NTK at initialization)

Fix  $\epsilon > 0$  and  $\delta \in (0, 1)$ . Suppose  $\sigma(z) = \max(0, z)$  and  $\min_{h \in [L]} d_h \geq \Omega(\frac{L^6}{\epsilon^4} \log(L/\delta))$ . 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:

$$\left| \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}')}{\partial \boldsymbol{\theta}} \right\rangle - \Theta^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L + 1)\epsilon.$$

**Improvements upon previous results (Jacot et al. 2018):**

- **Non-asymptotic bound** on the required layer widths, previous results are asymptotic, (i.e. all  $d_h \rightarrow \infty$ )
- Only requires  $\min_{h \in [L]} d_h$  to be sufficiently large, which is the weakest notion of limit.



# Equivalence wide neural net and kernel regression

- Can further incorporate the training process and show the equivalence between a fully-trained sufficiently wide neural net and the kernel regression solution

## Equivalence wide neural net and kernel regression

- Can further incorporate the training process and show the equivalence between a fully-trained sufficiently wide neural net and the kernel regression solution
- Final output of the neural network be  $f_{nn}(\boldsymbol{\theta}, \mathbf{x}) = \kappa f(\boldsymbol{\theta}, \mathbf{x})$ . and  $f_{nn}(\mathbf{x}_{te}) = \lim_{t \rightarrow \infty} f_{nn}(\boldsymbol{\theta}(t), \mathbf{x}_{te})$  be the prediction of the neural network at the end of training

## Equivalence wide neural net and kernel regression

Theorem (2) Equivalence between trained net and kernel regression)

Suppose  $\sigma(z) = \max(0, z)$ ,  $1/\kappa = \text{poly}(1/\epsilon, \log(n/\delta))$  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

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

## Equivalence wide neural net and kernel regression

Theorem (2) Equivalence between trained net and kernel regression)

Suppose  $\sigma(z) = \max(0, z)$ ,  $1/\kappa = \text{poly}(1/\epsilon, \log(n/\delta))$  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

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

- **Remark:** Result can be extended to *(exponentially many)* *finite testing points* using a union bound.

## Proof Idea Theorem 2

- First use a generic argument to show that the **perturbation on the prediction** can be reduced to the **perturbation on kernel value** at the **initialization** and during **training**.

## Proof Idea Theorem 2

- First use a generic argument to show that the **perturbation on the prediction** can be reduced to the **perturbation on kernel value** at the **initialization** and during **training**.
- **Initialization:** Theorem 1 guarantees a small perturbation on kernel value.
- **Training:** use high level proof idea from Du et al. 2018 to reduce the perturbation on the kernel value to the **perturbation on the gradient of each prediction with respect to weight matrices**. Then we adopt technical lemmas from Allen-Zhu et al. 2018 to obtain bounds on the perturbation of the gradient.

# 3. Proofs

## NTK Derivation

- **Goal:** Compute the limit of  $\left\langle \frac{\partial f(\theta, \mathbf{x})}{\partial \theta}, \frac{\partial f(\theta, \mathbf{x}')}{\partial \theta} \right\rangle$  at random initialization in the infinite width limit.
- Recall:  $\mathbf{f}^{(h)}(\mathbf{x}) = \mathbf{W}^{(h)} \mathbf{g}^{(h-1)}(\mathbf{x}) \in \mathbb{R}^{d_h}$



# NTK Derivation

- **Goal:** Compute the limit of  $\left\langle \frac{\partial f(\theta, \mathbf{x})}{\partial \theta}, \frac{\partial f(\theta, \mathbf{x}')}{\partial \theta} \right\rangle$  at random initialization in the infinite width limit.
- Recall:  $\mathbf{f}^{(h)}(\mathbf{x}) = \mathbf{W}^{(h)} \mathbf{g}^{(h-1)}(\mathbf{x}) \in \mathbb{R}^{d_h}$
- Partial derivative w.r.t. particular weight  $\mathbf{W}^{(h)}$ :

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

where

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

$$\mathbf{D}^{(h)}(\mathbf{x}) = \text{diag} \left( \dot{\sigma} \left( \mathbf{f}^{(h)}(\mathbf{x}) \right) \right) \in \mathbb{R}^{d_h \times d_h}, \quad h = 1, \dots, L.$$

## NTK Derivation

Then, for any  $h \in [L + 1]$ , we can compute

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

- Using CLT,

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

- Inductively, can show that

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

## Proof Theorem 2

- Recall: Final output of the neural network be  $f_{nn}(\boldsymbol{\theta}, \mathbf{x}) = \kappa f(\boldsymbol{\theta}, \mathbf{x})$ . and  $f_{nn}(\mathbf{x}_{te}) = \lim_{t \rightarrow \infty} f_{nn}(\boldsymbol{\theta}(t), \mathbf{x}_{te})$  be the prediction of the neural network at the end of training.
- We define  $\ker_{\mathbf{t}}(\mathbf{x}_{te}, \mathbf{X}) \in \mathbb{R}^n$  as

$$[\ker_{\mathbf{t}}(\mathbf{x}_{te}, \mathbf{X})]_i = \left\langle \frac{\partial f(\boldsymbol{\theta}(\mathbf{t}), \mathbf{x}_{te})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}(\mathbf{t}), \mathbf{x}_i)}{\partial \boldsymbol{\theta}} \right\rangle$$

## Proof Theorem 2

- Recall: Final output of the neural network be  $f_{nn}(\boldsymbol{\theta}, \mathbf{x}) = \kappa f(\boldsymbol{\theta}, \mathbf{x})$ . and  $f_{nn}(\mathbf{x}_{te}) = \lim_{t \rightarrow \infty} f_{nn}(\boldsymbol{\theta}(t), \mathbf{x}_{te})$  be the prediction of the neural network at the end of training.
- We define  $\ker_t(\mathbf{x}_{te}, \mathbf{X}) \in \mathbb{R}^n$  as

$$[\ker_t(\mathbf{x}_{te}, \mathbf{X})]_i = \left\langle \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_{te})}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}(t), \mathbf{x}_i)}{\partial \boldsymbol{\theta}} \right\rangle$$

- Let  $\ker_{ntk}(\mathbf{x}_{te}, \mathbf{X}) \in \mathbb{R}^n$  be the kernel values between  $\mathbf{x}_{te}$  and each training data (i.e.,  $[\ker_{ntk}(\mathbf{x}_{te}, \mathbf{X})]_i = \Theta^{(L)}(\mathbf{x}_{te}, \mathbf{x}_i)$ ). Then, the prediction of kernel regression is

$$f_{ntk}(\mathbf{x}_{te}) = \ker_{ntk}(\mathbf{x}_{te}, \mathbf{X})^\top \underbrace{(\mathbf{H}^*)^{-1}}_{[\mathbf{H}^*]_{i,j} = \Theta^{(L)}(\mathbf{x}_i, \mathbf{x}_j)} \mathbf{y}. \quad (9)$$

## Proof Theorem 2

- First prove a lemma to reduce the prediction perturbation bound to the kernel perturbation bound.

Lemma (Kernel Value Perturbation  $\Rightarrow$  Prediction Perturbation)

Fix  $\epsilon_H \leq \frac{1}{2}\lambda_0$ . Suppose  $|f_{nn}(\boldsymbol{\theta}(0), \mathbf{x}_i)| \leq \epsilon_{init}$  for  $i = 1, \dots, n$  and  $|f_{nn}(\boldsymbol{\theta}(0), \mathbf{x}_{te})| \leq \epsilon_{init}$  and  $\|\mathbf{u}_{nn}(0) - \mathbf{y}\|_2 = O(\sqrt{n})$ .

Furthermore, if for all  $t \geq 0$

$\|\ker_{ntk}(\mathbf{x}_{te}, \mathbf{X}) - \ker_t(\mathbf{x}_{te}, \mathbf{X})\|_2 \leq \epsilon_{test}$  and

$\|\mathbf{H}^* - \mathbf{H}(t)\|_2 \leq \epsilon_H$ , then we have

$$|f_{ntk}(\mathbf{x}_{te}) - f_{nn}(\mathbf{x}_{te})| \leq O\left(\epsilon_{init} + \frac{\sqrt{n}}{\lambda_0} \epsilon_{test} + \frac{\sqrt{n}}{\lambda_0^2} \log\left(\frac{n}{\epsilon_H \lambda_0 \kappa}\right) \epsilon_H\right).$$

## Proof Theorem 2

- Use  $\beta(t)$  to denote this parameter at time  $t$  **trained by gradient flow** and  $f_{ntk}(\mathbf{x}_{te}, \beta(t))$  be the predictor for  $\mathbf{x}_{te}$  at time  $t$ .
- Then

$$f_{ntk}(\mathbf{x}_{te}) = \int_{t=0}^{\infty} \frac{df_{ntk}(\beta(t), \mathbf{x}_{te})}{dt} dt$$

- Now we take a closer look at the time derivative:

$$\frac{df_{ntk}(\beta(t), \mathbf{x}_{te})}{dt} = -\frac{\kappa^2}{n} \text{ker}_{ntk}(\mathbf{x}_{te}, \mathbf{X})^\top (\mathbf{u}_{ntk}(t) - \mathbf{y})$$

where  $u_{ntk,i}(t) = f_{ntk}(\beta(t), \mathbf{x}_i)$  and  $\mathbf{u}_{ntk}(t) \in \mathbb{R}^n$  with  $[\mathbf{u}_{ntk}(t)]_i = u_{ntk,i}(t)$ .

## Proof Theorem 2

- Similarly, for the NN predictor:

$$\frac{df_{nn}(\boldsymbol{\theta}(t), \mathbf{x}_{te})}{dt} = -\frac{\kappa^2}{n} \text{ker}_t(\mathbf{x}_{te}, \mathbf{X})^\top (\mathbf{u}_{nn}(t) - \mathbf{y})$$

## Proof Theorem 2

- Similarly, for the NN predictor:

$$\frac{df_{nn}(\boldsymbol{\theta}(t), \mathbf{x}_{te})}{dt} = -\frac{\kappa^2}{n} \text{ker}_t(\mathbf{x}_{te}, \mathbf{X})^\top (\mathbf{u}_{nn}(t) - \mathbf{y})$$

- Next, we analyze the difference between the NN predictor and NTK predictor via this integral form

$$\begin{aligned} & |f_{nn}(\mathbf{x}_{te}) - f_{ntk}(\mathbf{x}_{te})| \\ & \leq \epsilon_{init} + \kappa^2 \epsilon_{test} \int_{t=0}^{\infty} \|\mathbf{u}_{nn}(t) - \mathbf{y}\|_2 dt \\ & \quad + \kappa^2 \max_{0 \leq t \leq \infty} \|\text{ker}_{ntk}(\mathbf{x}_{te}, \mathbf{X})\|_2 \int_{t=0}^{\infty} \|\mathbf{u}_{nn}(t) - \mathbf{u}_{ntk}(t)\|_2 dt \end{aligned}$$



## Proof Theorem 2

- Next, we observe that  $\mathbf{u}_{nn}(t) \rightarrow \mathbf{y}$  and  $\mathbf{u}_{ntk}(t) \rightarrow \mathbf{y}$  with linear convergence rate (Allen-Zhu et al. 2018), thus:

$$\|\mathbf{u}_{nn}(t) - \mathbf{y}\|_2 \leq \exp\left(-\frac{\kappa^2}{2}\lambda_0 t\right) \underbrace{\|\mathbf{u}_{nn}(0) - \mathbf{y}\|_2}_{O(\sqrt{n})}$$

.

## Concluding the proof

By previous Lemma, the problem now reduces to

- i) Choose  $\kappa$  small enough to make  $\epsilon_{init} = O(\epsilon)$
- ii) Show when the width is large enough then  $\epsilon_H$  and  $\epsilon_{test}$  are both  $O(\epsilon)$ .

## Concluding the proof

By previous Lemma, the problem now reduces to

- i) Choose  $\kappa$  small enough to make  $\epsilon_{init} = O(\epsilon)$
- ii) Show when the width is large enough then  $\epsilon_H$  and  $\epsilon_{test}$  are both  $O(\epsilon)$ .

How?

- i) Use result by Daniely et al., 2016 to choose  $\kappa = O\left(\frac{\epsilon}{\log(n/\delta)}\right)$  to make  $\epsilon_{init} = O(\epsilon)$  with probability  $1 - \delta$ .
- ii) Use Theorem 1 and additional lemma that bounds kernel perturbation

# Proof Theorem 1

- Infinite limit of the blue and red terms can be made more precise.

## Theorem

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.  $1 - \delta$ ,  $\forall 0 \leq h \leq L$ ,  
 $\forall (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{(\mathbf{x}, \mathbf{x}), (\mathbf{x}, \mathbf{x}'), (\mathbf{x}', \mathbf{x}')\}$ ,

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

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

# 4. Experiments

## Experiments on CIFAR-10 dataset

- Two CNN architectures: with and without global average pooling (GAP)
- Also test CNTKs with only 2,000 training data

## Experiments on CIFAR-10 dataset

- Two CNN architectures: with and without global average pooling (GAP)
- Also test CNTKs with only 2,000 training data

L	CNN-V	CNTK-V	CNTK-V-2K	CNN-GAP	CNTK-GAP	CNTK-GAP
3	59.97%	64.47%	40.94%	63.81%	70.47%	49.71%
4	60.20%	65.52%	42.54%	80.93%	75.93%	51.06%
6	64.11%	66.03%	43.43%	83.75%	76.73%	51.73%
11	69.48%	65.90%	43.42%	82.92%	<b>77.43%</b>	51.92%
21	75.57%	64.09%	42.53%	83.30%	77.08%	52.22%

- CNN-V represents vanilla CNN and CNTK-V represents the kernel corresponding to CNN-V.
- CNN-GAP represents CNN with GAP and CNTK-GAP represents the kernel corresponding to CNN-GAP.

## Experiments on CIFAR-10 dataset

- CNTKs are very powerful kernel, +10% compared to previous kernels
- For both CNN and CNTK, depth can affect the classification accuracy
- The global average pooling operation can significantly increase the classification accuracy by 8% – 10% for both CNN and CNT
- **Still 5% - 6% performance gap between CNTKs & CNNs**