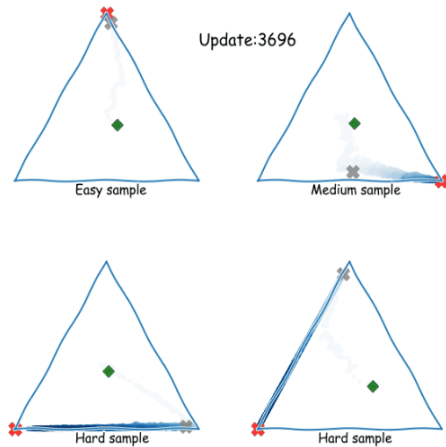# In Defence of (Empirical) Neural Tangent Kernels

**Danica J. Sutherland** (she)
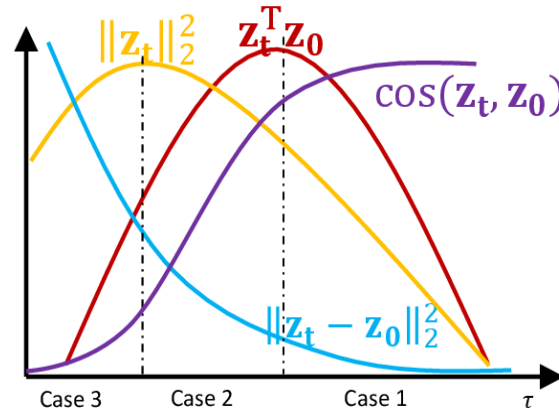
University of British Columbia (UBC) / Alberta Machine Intelligence Institute (Amii)

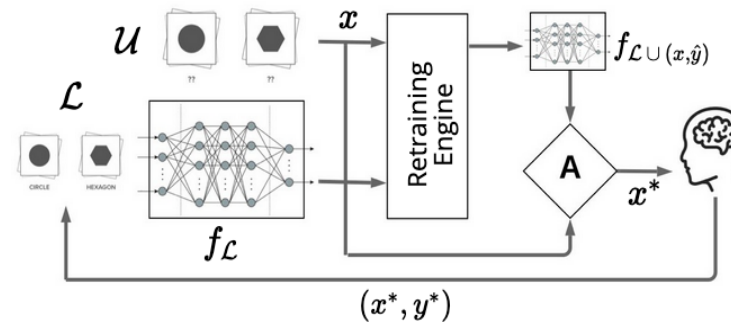"Zig-zagging" [ICLR-22]  Finetuning [ICLR-23]  Active learning [NeurIPS-22]  Pseudo-NTK [new!]



Yi Ren

Shangmin Guo

Yi Ren

Shangmin Guo
Wonho Bae

Mohamad Amin Mohamadi
Wonho Bae

Mohamad Amin Mohamadi
Wonho Bae

MSR Montréal / Mila - March 21, 2023

1

# One path to NTKs

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

# One path to NTKs

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t})}_{k \times p} \underbrace{(w_{t+1} - w_t)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{\left(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t}\right)}_{k \times p} \underbrace{\left(w_{t+1} - w_t\right)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

$$= \left(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}\right)\left(-\eta \nabla_{\mathbf{w}} \ell_{y_i}(f(x_i))|_{\mathbf{w}_t}\right)^{\mathsf{T}} + \mathcal{O}(\eta^2)$$

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{\left(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t}\right)}_{k \times p} \underbrace{\left(w_{t+1} - w_t\right)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

$$= \left(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}\right) \left(-\eta \nabla_{\mathbf{w}} \ell_{y_i}(f(x_i))|_{\mathbf{w}_t}\right)^{\mathsf{T}} + \mathcal{O}(\eta^2)$$

$$= -\eta \underbrace{\left(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}\right)}_{k \times p} \underbrace{\left(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t}\right)^{\mathsf{T}}}_{p \times k} \underbrace{\ell'_{y_i}(f_t(x_i))}_{k \times 1} + \mathcal{O}(\eta^2)$$

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{(\nabla_\mathbf{w} f(x_i)|_{\mathbf{w}_t})}_{k \times p} \underbrace{(w_{t+1} - w_t)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

$$= (\nabla_\mathbf{w} f(\tilde{x})|_{\mathbf{w}_t}) \left(-\eta \nabla_\mathbf{w} \ell_{y_i}(f(x_i))|_{\mathbf{w}_t}\right)^\mathsf{T} + \mathcal{O}(\eta^2)$$

$$= -\eta \underbrace{(\nabla_\mathbf{w} f(\tilde{x})|_{\mathbf{w}_t})}_{k \times p} \underbrace{(\nabla_\mathbf{w} f(x_i)|_{\mathbf{w}_t})^\mathsf{T}}_{p \times k} \underbrace{\ell'_{y_i}(f_t(x_i))}_{k \times 1} + \mathcal{O}(\eta^2)$$

$$= -\eta \, \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \, \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

- Defined $\mathrm{eNTK}_\mathbf{w}(\tilde{x}, x_i) = (\nabla_\mathbf{w} f(\tilde{x}; \mathbf{w})) (\nabla_\mathbf{w} f(x_i; \mathbf{w}))^\mathsf{T} \in \mathbb{R}^{k \times k}$

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t})}_{k \times p} \underbrace{(w_{t+1} - w_t)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

$$= (\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}) \left(-\eta \nabla_{\mathbf{w}} \ell_{y_i}(f(x_i))|_{\mathbf{w}_t}\right)^{\mathsf{T}} + \mathcal{O}(\eta^2)$$

$$= -\eta \underbrace{(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t})}_{k \times p} \underbrace{(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t})^{\mathsf{T}}}_{p \times k} \underbrace{\ell'_{y_i}(f_t(x_i))}_{k \times 1} + \mathcal{O}(\eta^2)$$

$$= -\eta \, \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \, \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

- Defined $\mathrm{eNTK}_{\mathbf{w}}(\tilde{x}, x_i) = (\nabla_{\mathbf{w}} f(\tilde{x}; \mathbf{w})) (\nabla_{\mathbf{w}} f(x_i; \mathbf{w}))^{\mathsf{T}} \in \mathbb{R}^{k \times k}$
  - $(x_i, y_i)$ step barely changes $\tilde{x}$ prediction if $\mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)$ is small

# One path to NTKs (Taylor's version)

- "Learning path" of a model's predictions: $f_t(\tilde{x})$ for some fixed $\tilde{x}$ as params $\mathbf{w}_t$ change

- Let's start with "plain" SGD on $\frac{1}{N} \sum_{i=1}^{N} \ell_{y_i}(f(x_i))$:

$$\underbrace{f_{t+1}(\tilde{x}) - f_t(\tilde{x})}_{k \times 1} = \underbrace{\left(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t}\right)}_{k \times p} \underbrace{\left(w_{t+1} - w_t\right)}_{p \times 1} + \mathcal{O}\left(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2\right)$$

$$= \left(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}\right) \left(-\eta \nabla_{\mathbf{w}} \ell_{y_i}(f(x_i))|_{\mathbf{w}_t}\right)^{\mathsf{T}} + \mathcal{O}(\eta^2)$$

$$= -\eta \underbrace{\left(\nabla_{\mathbf{w}} f(\tilde{x})|_{\mathbf{w}_t}\right)}_{k \times p} \underbrace{\left(\nabla_{\mathbf{w}} f(x_i)|_{\mathbf{w}_t}\right)^{\mathsf{T}}}_{p \times k} \underbrace{\ell'_{y_i}(f_t(x_i))}_{k \times 1} + \mathcal{O}(\eta^2)$$

$$= -\eta \, \text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \, \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

- Defined $\text{eNTK}_{\mathbf{w}}(\tilde{x}, x_i) = \left(\nabla_{\mathbf{w}} f(\tilde{x}; \mathbf{w})\right) \left(\nabla_{\mathbf{w}} f(x_i; \mathbf{w})\right)^{\mathsf{T}} \in \mathbb{R}^{k \times k}$
  - $(x_i, y_i)$ step barely changes $\tilde{x}$ prediction if $\text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)$ is small

- $\ell'_y(\hat{y}) = \hat{y} - y$ for square loss, $\hat{y}_y - \log \sum_{j=1}^{k} \exp(\hat{y}_j)$ for cross-entropy

- Full-batch GD:

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N} \sum_{i=1}^{N} \text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

- Full-batch GD: "stacking things up",

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N}\sum_{i=1}^{N}\mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)\ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

$$= -\frac{\eta}{N}\underbrace{\mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, \mathbf{X})}_{k \times kN}\underbrace{L'_{\mathbf{y}}(f_t(\mathbf{X}))}_{kN \times 1} + \mathcal{O}(\eta^2)$$

- Full-batch GD: "stacking things up",

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N} \sum_{i=1}^{N} \text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)\ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

$$= -\frac{\eta}{N} \underbrace{\text{eNTK}_{\mathbf{w}_t}(\tilde{x}, \mathbf{X})}_{k \times kN} \underbrace{L'_{\mathbf{y}}(f_t(\mathbf{X}))}_{kN \times 1} + \mathcal{O}(\eta^2)$$

- Observation I: If $f$ is "wide enough" with any usual architecture+init* [Yang+Litwin 2021], $\text{eNTK}(\cdot, \mathbf{X})$ is roughly constant through training

- Full-batch GD: "stacking things up",

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N} \sum_{i=1}^{N} \text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)\ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

$$= -\frac{\eta}{N} \underbrace{\text{eNTK}_{\mathbf{w}_t}(\tilde{x}, \mathbf{X})}_{k \times kN} \underbrace{L'_{\mathbf{y}}(f_t(\mathbf{X}))}_{kN \times 1} + \mathcal{O}(\eta^2)$$

- Observation I: If $f$ is "wide enough" with any usual architecture+init* [Yang+Litwin 2021], $\text{eNTK}(\cdot, \mathbf{X})$ is roughly constant through training
  - For square loss, $L'_{y_i}(f_t(\mathbf{X})) = f_t(\mathbf{X}) - y_i$: dynamics agree with kernel regression!

- Full-batch GD: "stacking things up",

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N} \sum_{i=1}^{N} \text{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

$$= -\frac{\eta}{N} \underbrace{\text{eNTK}_{\mathbf{w}_t}(\tilde{x}, \mathbf{X})}_{k \times kN} \underbrace{L'_{\mathbf{y}}(f_t(\mathbf{X}))}_{kN \times 1} + \mathcal{O}(\eta^2)$$

- Observation I: If $f$ is "wide enough" with any usual architecture+init* [Yang+Litwin 2021], $\text{eNTK}(\cdot, \mathbf{X})$ is roughly constant through training
  - For square loss, $L'_{y_i}(f_t(\mathbf{X})) = f_t(\mathbf{X}) - y_i$: dynamics agree with kernel regression!
  - $f_t(\tilde{x}) \xrightarrow{t \to \infty} \text{eNTK}_{\mathbf{w}_0}(\tilde{x}, \mathbf{X}) \, \text{eNTK}_{\mathbf{w}_0}(\mathbf{X}, \mathbf{X})^{-1}(\mathbf{y} - f_0(\mathbf{X})) + f_0(\tilde{x})$

- Full-batch GD: "stacking things up",

$$f_{t+1}(\tilde{x}) - f_t(\tilde{x}) = -\frac{\eta}{N} \sum_{i=1}^{N} \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \ell'_{y_i}(f_t(x_i)) + \mathcal{O}(\eta^2)$$

$$= -\frac{\eta}{N} \underbrace{\mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, \mathbf{X})}_{k \times kN} \underbrace{L'_{\mathbf{y}}(f_t(\mathbf{X}))}_{kN \times 1} + \mathcal{O}(\eta^2)$$

- Observation I: If $f$ is "wide enough" with any usual architecture+init* [Yang+Litwin 2021], $\mathrm{eNTK}(\cdot, \mathbf{X})$ is roughly constant through training
  - For square loss, $L'_{y_i}(f_t(\mathbf{X})) = f_t(\mathbf{X}) - y_i$: dynamics agree with kernel regression!
  - $f_t(\tilde{x}) \xrightarrow{t \to \infty} \mathrm{eNTK}_{\mathbf{w}_0}(\tilde{x}, \mathbf{X}) \, \mathrm{eNTK}_{\mathbf{w}_0}(\mathbf{X}, \mathbf{X})^{-1} (\mathbf{y} - f_0(\mathbf{X})) + f_0(\tilde{x})$

- Observation II: As $f$ becomes "infinitely wide" with any usual architecture+init* [Yang 2019], $\mathrm{eNTK}_{\mathbf{w}_0}(x_1, x_2) \xrightarrow{a.s.} \mathrm{NTK}(x_1, x_2)$, independent of the random $\mathbf{w}_0$

# Infinite NTKs are great

- Infinitely-wide neural networks have very simple behaviour!
  - No need to worry about bad local minima, optimization complications, …

# **Infinite NTKs are great**

- Infinitely-wide neural networks have very simple behaviour!
  - No need to worry about bad local minima, optimization complications, …
  - Understanding "implicit bias" of wide nets ≈ understanding NTK norm of functions

# Infinite NTKs are great

- Infinitely-wide neural networks have <span style="color:green">very simple behaviour</span>!
  - No need to worry about bad local minima, optimization complications, ...

  - Understanding "implicit bias" of wide nets $\approx$ understanding NTK norm of functions

- Can compute $\mathrm{NTK}$ exactly for many architectures
  - github.com/google/neural-tangents

# Infinite NTKs are great

- Infinitely-wide neural networks have very simple behaviour!
  - No need to worry about bad local minima, optimization complications, …

  - Understanding "implicit bias" of wide nets ≈ understanding NTK norm of functions

- Can compute $\mathbf{NTK}$ exactly for many architectures
  - github.com/google/neural-tangents

- A great kernel for many kernel methods!
  - Using in SVMs was then-best overall method across many small-data tasks [Arora+ 2020]

# Infinite NTKs are great

- Infinitely-wide neural networks have <span style="color:green">very simple behaviour</span>!
  - No need to worry about bad local minima, optimization complications, ...

  - Understanding "implicit bias" of wide nets $\approx$ understanding NTK norm of functions

- Can compute $\mathrm{NTK}$ exactly for many architectures
  - `github.com/google/neural-tangents`

- A great kernel for many kernel methods!
  - Using in SVMs was then-best overall method across many small-data tasks [Arora+ 2020]

  - Good results in statistical testing [Jia+ 2021], dataset distillation [Nguyen+ 2021], clustering for active learning batch queries [Holzmüller+ 2022], ...

# But (infinite) NTKs aren't "the answer"

- Computational expense:
  - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
    - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!

# But (infinite) NTKs aren't "the answer"

- Computational expense:
    - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
        - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!

        - ILSVRC2012 has $n \approx 1\,200\,000$, $k = 1\,000$: 11.5 *million* terabytes (exabytes)

# But (infinite) NTKs aren't "the answer"

- Computational expense:
  - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
    - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!
    - ILSVRC2012 has $n \approx 1\,200\,000$, $k = 1\,000$: 11.5 *million* terabytes (exabytes)
  - For deep/complex models (especially CNNs), each pair very slow / memory-intensive

# But (infinite) NTKs aren't "the answer"

- Computational expense:
  - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
    - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!

    - ILSVRC2012 has $n \approx 1\,200\,000$, $k = 1\,000$: 11.5 *million* terabytes (exabytes)

  - For deep/complex models (especially CNNs), each pair very slow / memory-intensive

- Practical performance:
  - Typically performs worse than GD for "non-small-data" tasks (MNIST and up)

# But (infinite) NTKs aren't "the answer"

- Computational expense:
  - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
    - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!
    - ILSVRC2012 has $n \approx 1\,200\,000$, $k = 1\,000$: 11.5 *million* terabytes (exabytes)
  - For deep/complex models (especially CNNs), each pair very slow / memory-intensive
- Practical performance:
  - Typically performs worse than GD for "non-small-data" tasks (MNIST and up)
- Theoretical limitations:
  - NTK "doesn't do feature learning":
    - **eNTK** stays $\approx$ constant
    - Internal activations in the networks don't change much [Chizat+ 2019] [Yang/Hu 2021]

# But (infinite) NTKs aren't "the answer"

- Computational expense:
  - Poor scaling for large-data problems: typically $n^2$ memory and $n^2$ to $n^3$ computation
    - CIFAR-10 has $n = 50\,000$, $k = 10$: an $nk \times nk$ matrix of `float64`s is 2 terabytes!
    - ILSVRC2012 has $n \approx 1\,200\,000$, $k = 1\,000$: 11.5 *million* terabytes (exabytes)
  - For deep/complex models (especially CNNs), each pair very slow / memory-intensive

- Practical performance:
  - Typically performs worse than GD for "non-small-data" tasks (MNIST and up)

- Theoretical limitations:
  - NTK "doesn't do feature learning":
    - **eNTK** stays $\approx$ constant
    - Internal activations in the networks don't change much [Chizat+ 2019] [Yang/Hu 2021]
  - We now know many problems where gradient descent on an NN $\gg$ *any* kernel method
    - Cases where GD error $\rightarrow 0$, any kernel is *barely* better than random [Malach+ 2021]

# What can we learn from empirical NTKs?

In this talk:

- As a theoretical-ish tool for local understanding:
    - Fine-grained explanation for early stopping in knowledge distillation
    - How you should fine-tune models
- As a practical tool for approximating "lookahead" in active learning
- Plus: efficiently approximating $\mathbf{eNTK}$s for large output dimensions $k$, with guarantees

# **What can we learn from empirical NTKs?**

In this talk:

- As a theoretical-ish tool for local understanding:
  - Fine-grained explanation for early stopping in knowledge distillation

  - How you should fine-tune models

- As a practical tool for approximating "lookahead" in active learning

- Plus: efficiently approximating $\mathbf{eNTK}$s for large output dimensions $k$, with guarantees

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)} \, \ell(f(x), y) = \mathbb{E}_x \, \mathbb{E}_{y|x} \, \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i), y_i)$$

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)} \ell(f(x), y) = \mathbb{E}_x \mathbb{E}_{y|x} \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i), y_i) = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{k} \mathbb{I}(y_i = c) \ell(f(x_i), c)$$

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)} \, \ell(f(x), y) = \mathbb{E}_x \, \mathbb{E}_{y|x} \, \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize ($\vec{\ell}(\hat{y}) \in \mathbb{R}^k$ is vector of losses for all possible labels)

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i), y_i) = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{k} \mathbb{I}(y_i = c) \ell(f(x_i), c) = \frac{1}{N} \sum_{i=1}^{n} e_{y_i} \cdot \vec{\ell}(f(x_i))$$

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)} \, \ell(f(x), y) = \mathbb{E}_x \, \mathbb{E}_{y|x} \, \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize ($\vec{\ell}(\hat{y}) \in \mathbb{R}^k$ is vector of losses for all possible labels)

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i), y_i) = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{k} \mathbb{I}(y_i = c) \ell(f(x_i), c) = \frac{1}{N} \sum_{i=1}^{n} e_{y_i} \cdot \vec{\ell}(f(x_i))$$
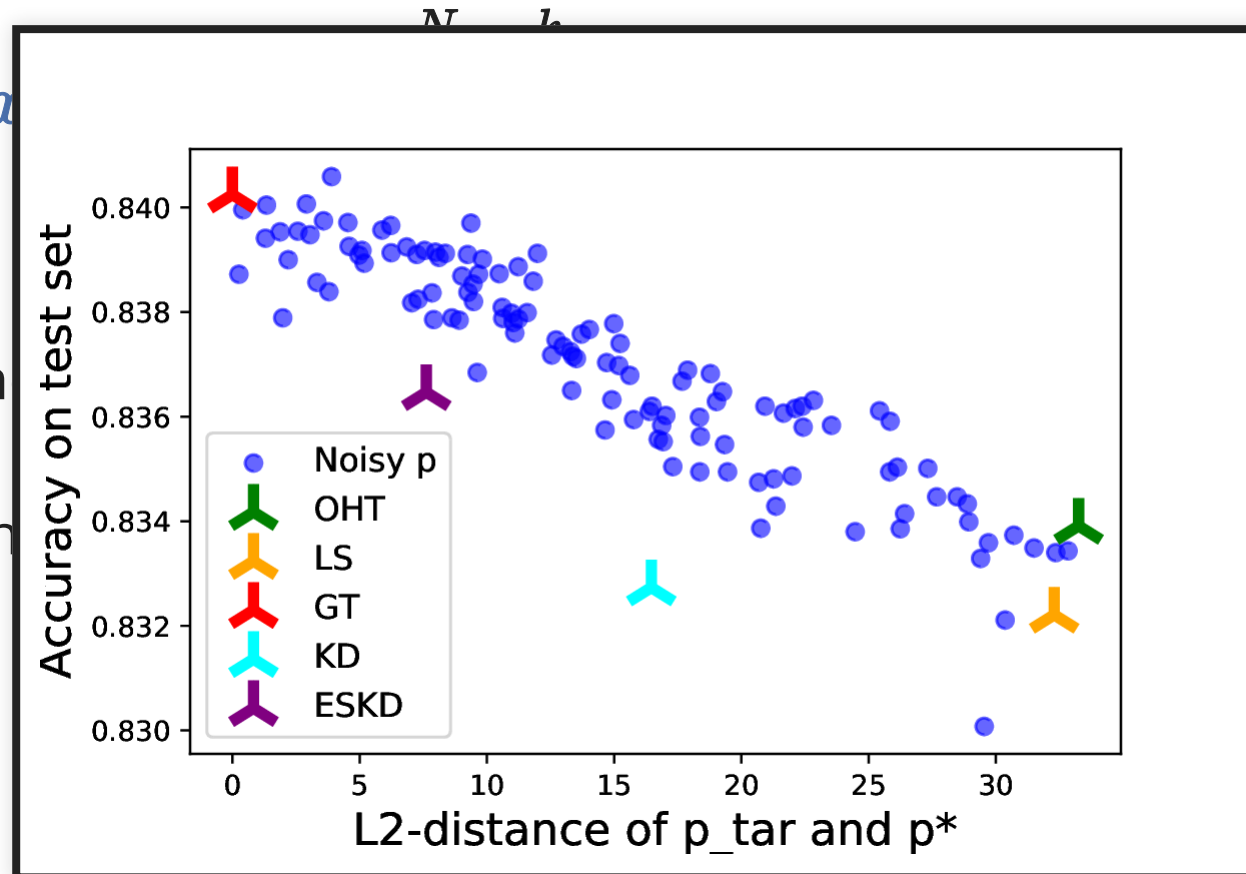
- Potentially better scheme: see $\{(x_i, p_i^{tar})\}$, minimize $L^{tar}(f) = \frac{1}{N} \sum_{i=1}^{N} p_i^{tar} \cdot \vec{\ell}(f(x_i))$

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)}\, \ell(f(x), y) = \mathbb{E}_x\, \mathbb{E}_{y|x}\, \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize ($\vec{\ell}(\hat{y}) \in \mathbb{R}^k$ is vector of losses for all possible labels)

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i), y_i) = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{k} \mathbb{I}(y_i = c)\ell(f(x_i), c) = \frac{1}{N} \sum_{i=1}^{n} e_{y_i} \cdot \vec{\ell}(f(x_i))$$

- Potentially better scheme: see $\{(x_i, p_i^{tar})\}$, minimize $L^{tar}(f) = \frac{1}{N} \sum_{i=1}^{N} p_i^{tar} \cdot \vec{\ell}(f(x_i))$

  - Can reduce variance if $p_i^{tar} \approx p_i^*$, the true conditional probabilities

# Better supervisory signal implies better learning

- Classification: target is $L_P(f) = \mathbb{E}_{(x,y)} \, \ell(f(x), y) = \mathbb{E}_x \, \mathbb{E}_{y|x} \, \ell(f(x), y)$

- Normally: see $\{(x_i, y_i)\}$, minimize ($\vec{\ell}(\hat{y}) \in \mathbb{R}^k$ is vector of losses for all possible labels)

$$L(f) = \frac{1}{N} \sum_{i=1}^{N} \ell(f(x \qquad \qquad \frac{1}{N} \sum_{i=1}^{n} e_{y_i} \cdot \vec{\ell}(f(x_i))$$

- Potentially better sch $\qquad \qquad \sum_{i=1}^{N} p_i^{tar} \cdot \vec{\ell}(f(x_i))$
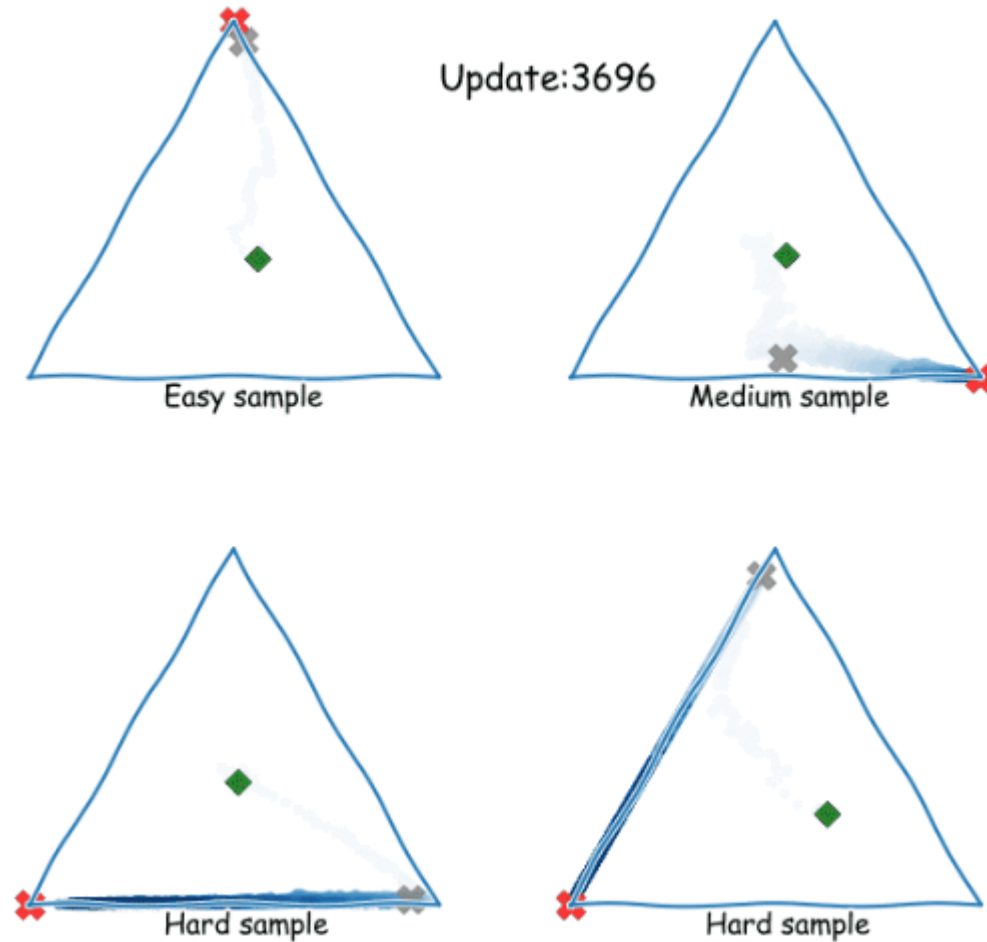
  - Can reduce varian

# Knowledge distillation

- Process:
  - Train a teacher $f^{teacher}$ on $\{(x_i, y_i)\}$ with standard ERM, $L(f)$

  - Train a student on $\{(x_i, f^{teacher}(x_i))\}$ with $L^{tar}$

- Usually $f^{student}$ is "smaller" than $f^{teacher}$

# Knowledge distillation

- Process:
  - Train a teacher $f^{teacher}$ on $\{(x_i, y_i)\}$ with standard ERM, $L(f)$
  - Train a student on $\{(x_i, f^{teacher}(x_i))\}$ with $L^{tar}$
- Usually $f^{student}$ is "smaller" than $f^{teacher}$
- But "self-distillation" (using the same architecture), often $f^{student}$ outperforms $f^{teacher}$!

# Knowledge distillation

- Process:
    - Train a teacher $f^{teacher}$ on $\{(x_i, y_i)\}$ with standard ERM, $L(f)$
    - Train a student on $\{(x_i, f^{teacher}(x_i))\}$ with $L^{tar}$

- Usually $f^{student}$ is "smaller" than $f^{teacher}$

- But "self-distillation" (using the same architecture), often $f^{student}$ outperforms $f^{teacher}$!

- One possible explanation: $f^{teacher}(x_i)$ is closer to $p_i^*$ than sampled $y_i$

# Knowledge distillation

- Process:
  - Train a teacher $f^{teacher}$ on $\{(x_i, y_i)\}$ with standard ERM, $L(f)$
  - Train a student on $\{(x_i, f^{teacher}(x_i))\}$ with $L^{tar}$
- Usually $f^{student}$ is "smaller" than $f^{teacher}$
- But "self-distillation" (using the same architecture), often $f^{student}$ outperforms $f^{teacher}$!
- One possible explanation: $f^{teacher}(x_i)$ is closer to $p_i^*$ than sampled $y_i$
- But why would that be?

# Zig-Zagging behaviour in learning



Plots of (three-way) probabilistic predictions: ✕ shows $p_i^*$, ✕ shows $y_i$

# eNTK explains it

- Let $q_t(\tilde{x}) = \mathrm{softmax}(f_t(\tilde{x})) \in \mathbb{R}^k$; for cross-entropy loss, one SGD step gives us

$$q_{t+1}(\tilde{x}) - q_t(\tilde{x}) = \eta \, \mathcal{A}_t(\tilde{x}) \, \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \, (p_i^{tar} - q_t(x_i)) + \mathcal{O}(\eta^2)$$
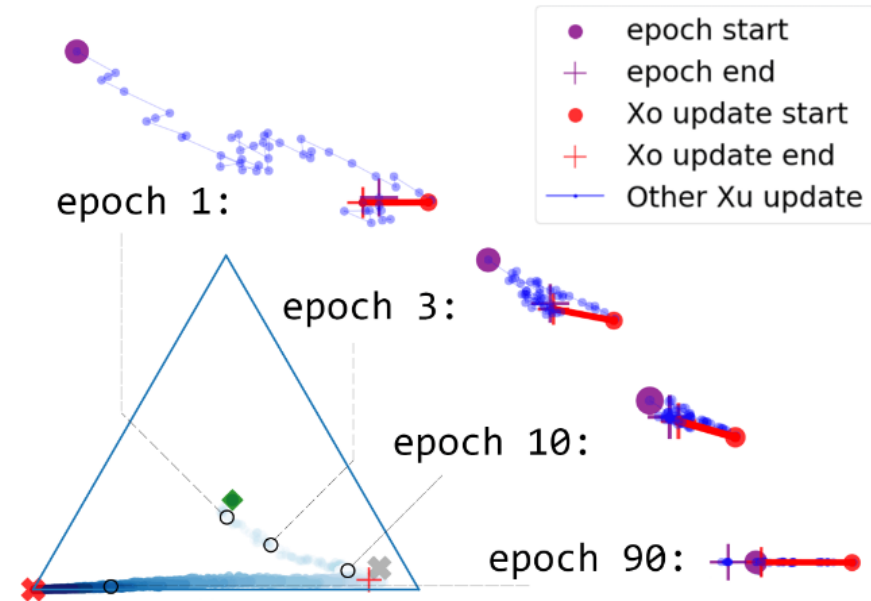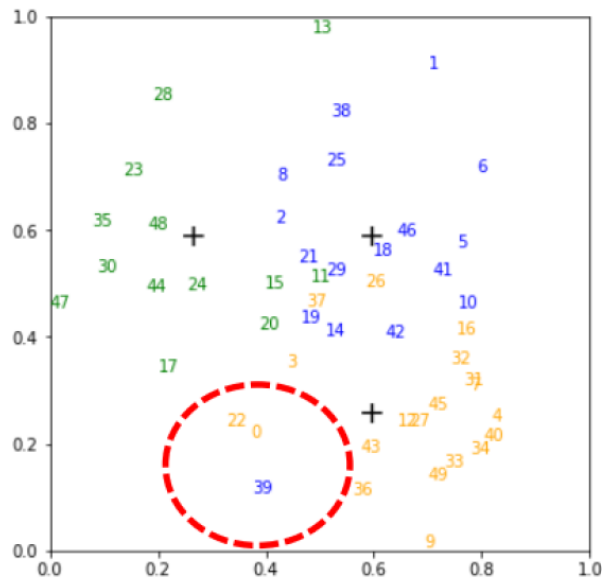
$$\mathcal{A}_t(\tilde{x}) = \mathrm{diag}(q_t(\tilde{x})) - q_t(\tilde{x}) q_t(\tilde{x})^{\mathsf{T}} \text{ is the covariance of a } \mathrm{Categorical}(q_t(\tilde{x}))$$
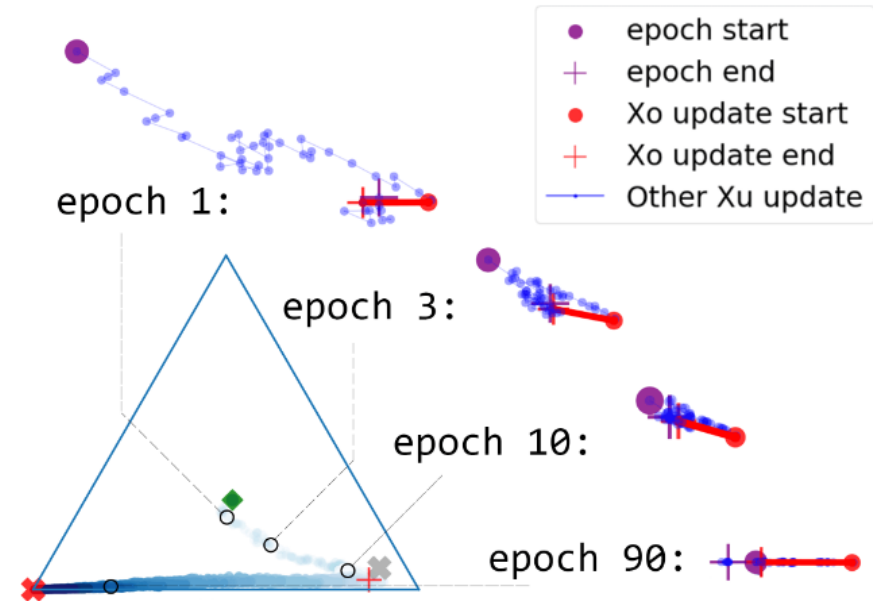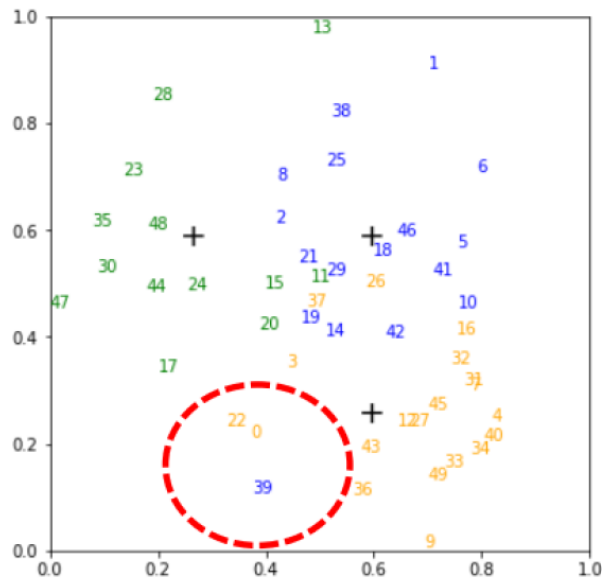
# eNTK explains it

- Let $q_t(\tilde{x}) = \mathrm{softmax}(f_t(\tilde{x})) \in \mathbb{R}^k$; for cross-entropy loss, one SGD step gives us

$$q_{t+1}(\tilde{x}) - q_t(\tilde{x}) = \eta \, \mathcal{A}_t(\tilde{x}) \, \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i) \left(p_i^{tar} - q_t(x_i)\right) + \mathcal{O}(\eta^2)$$

$$\mathcal{A}_t(\tilde{x}) = \mathrm{diag}(q_t(\tilde{x})) - q_t(\tilde{x})q_t(\tilde{x})^\mathsf{T} \text{ is the covariance of a } \mathrm{Categorical}(q_t(\tilde{x}))$$

# eNTK explains it

- Let $q_t(\tilde{x}) = \mathrm{softmax}(f_t(\tilde{x})) \in \mathbb{R}^k$; for cross-entropy loss, one SGD step gives us

$$q_{t+1}(\tilde{x}) - q_t(\tilde{x}) = \eta\, \mathcal{A}_t(\tilde{x})\, \mathrm{eNTK}_{\mathbf{w}_t}(\tilde{x}, x_i)\, (p_i^{tar} - q_t(x_i)) + \mathcal{O}(\eta^2)$$

$$\mathcal{A}_t(\tilde{x}) = \mathrm{diag}(q_t(\tilde{x})) - q_t(\tilde{x})q_t(\tilde{x})^{\mathsf{T}} \text{ is the covariance of a } \mathrm{Categorical}(q_t(\tilde{x}))$$
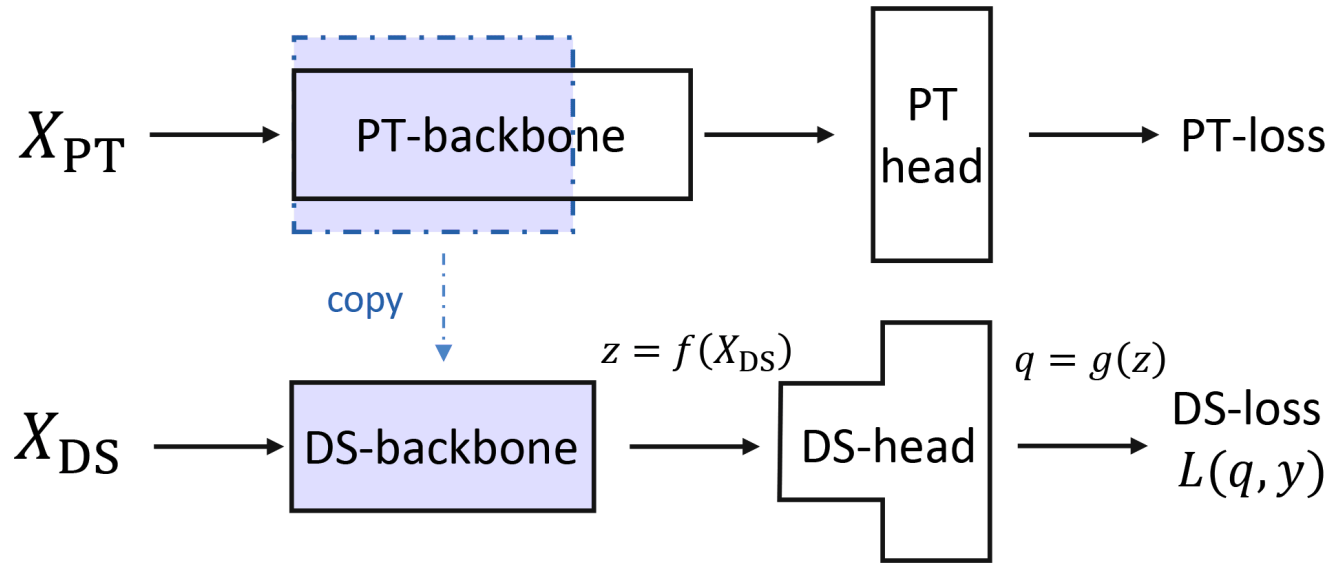


- Improves distillation (esp. with noisy labels) to take moving average of $q_t(x_i)$ as $p_i^{tar}$

# What can we learn from empirical NTKs?
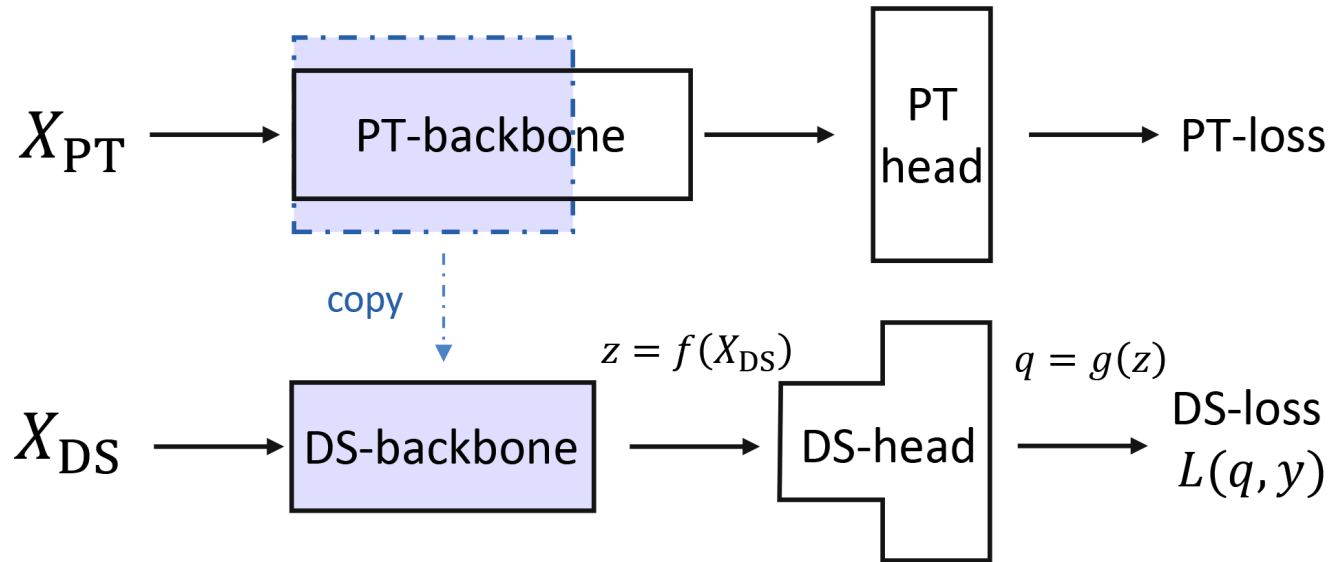
In this talk:

- As a theoretical-ish tool for local understanding:
  - Fine-grained explanation for early stopping in knowledge distillation
  - How you should fine-tune models

- As a practical tool for approximating "lookahead" in active learning

- Plus: efficiently approximating $\mathbf{eNTK}$s for large output dimensions $k$, with guarantees

# Fine-tuning



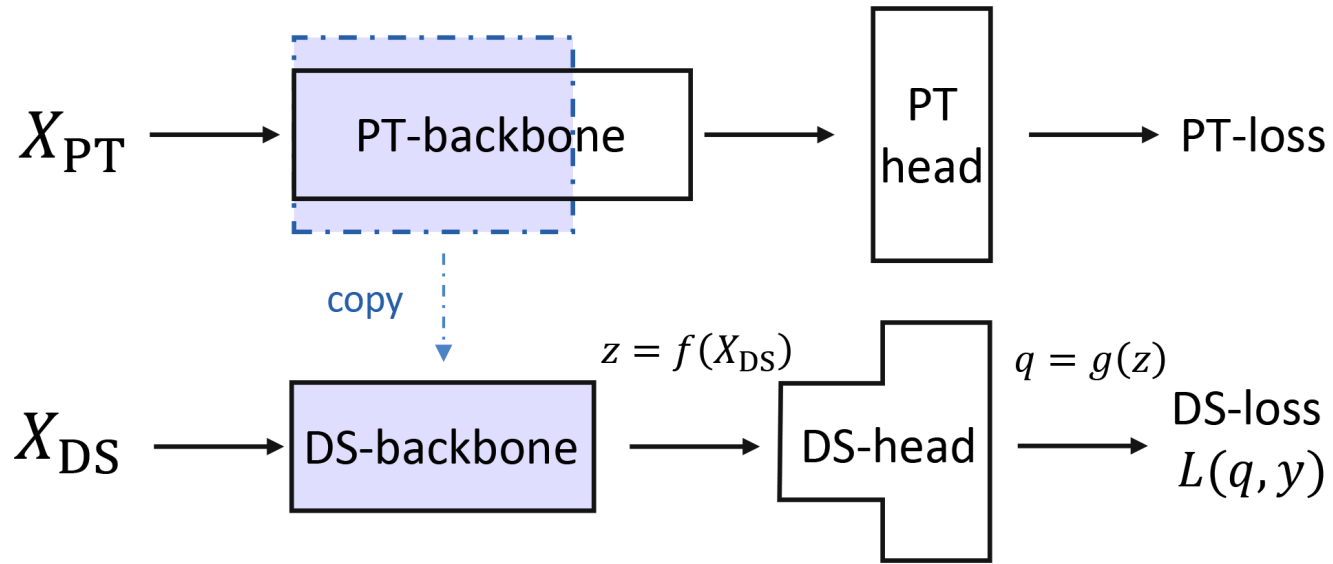- Pretrain, re-initialize a random head, then adapt to a downstream task. Two phases:
  - Head probing: only update the head $g(z)$
  - Fine-tuning: update head $g(z)$ and backbone $z = f(x)$ together

# Fine-tuning
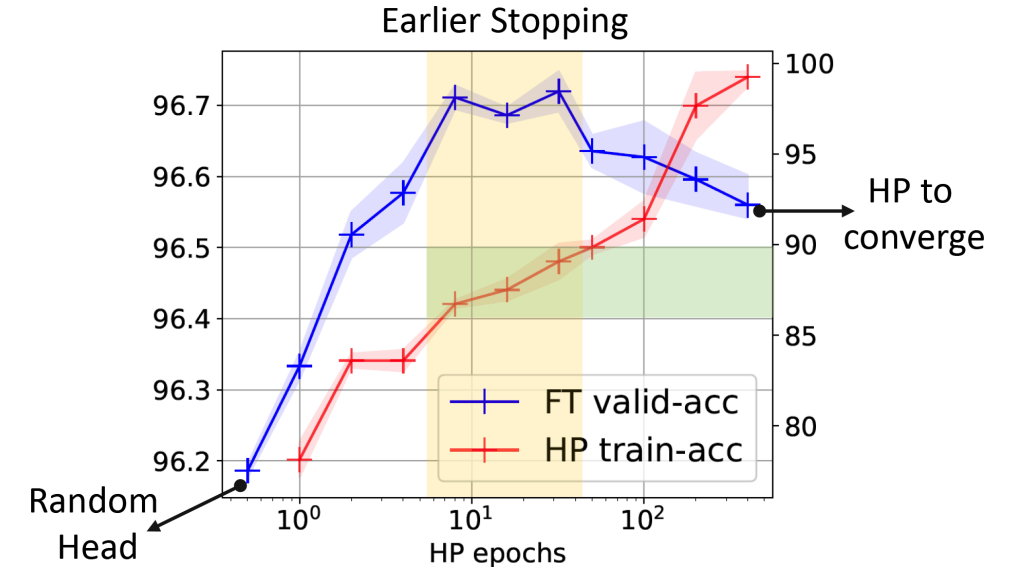


- Pretrain, re-initialize a random head, then adapt to a downstream task. Two phases:
  - Head probing: only update the head $g(z)$
  - Fine-tuning: update head $g(z)$ and backbone $z = f(x)$ together
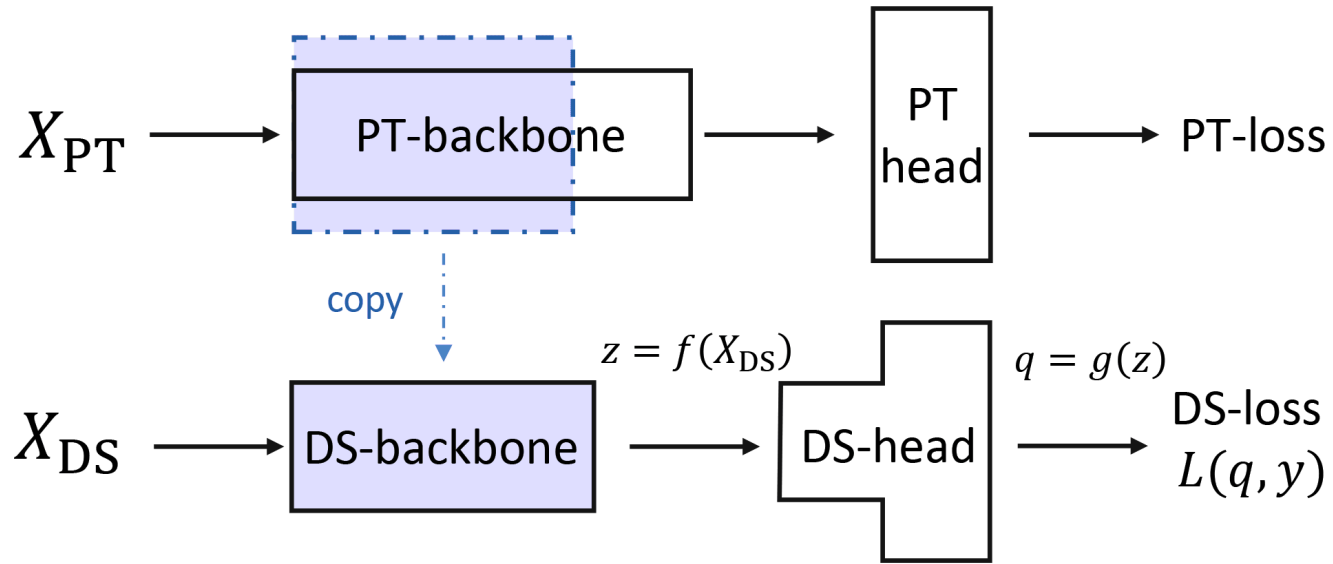- If we only fine-tune: noise from random head might break our features!

# Fine-tuning



- Pretrain, re-initialize a random head, then adapt to a downstream task. Two phases:
  - Head probing: only update the head $g(z)$

  - Fine-tuning: update head $g(z)$ and backbone $z = f(x)$ together

- If we only fine-tune: noise from random head might break our features!

- If we head-probe to convergence: might already fit training data and not change features!

# Fine-tuning



- Pretrain, re-initialize a random head, then adapt to a downstream task. Two phases:
  - Head probing: only update the head $g(z)$

  - Fine-tuning: update head $g(z)$ and backbone $z = f(x)$ together

- If we only fine-tune: noise from random head might break our features!

- If we head-probe to convergence: might already fit training data and not change features!

# How much do we change our features?

- Same kind of decomposition with backbone features $z = f(x)$, head $q = \mathrm{softmax}(g(z))$:

$$z_{t+1}(\tilde{x}) - z_t(\tilde{x}) = \frac{\eta}{N} \sum_{i=1}^{N} \underbrace{\mathrm{eNTK}^f_{\mathbf{w}_t}(\tilde{x}, x_i)}_{\text{eNTK of backbone}} \underbrace{(\nabla_z q_t(x_i))^\mathsf{T}}_{\text{direction of head}} \underbrace{(e_{y_i} - q_t(x_i))}_{\text{``energy''}} + \mathcal{O}(\eta^2)$$

# How much do we change our features?

- Same kind of decomposition with backbone features $z = f(x)$, head $q = \mathrm{softmax}(g(z))$:

$$z_{t+1}(\tilde{x}) - z_t(\tilde{x}) = \frac{\eta}{N} \sum_{i=1}^{N} \underbrace{\mathrm{eNTK}_{\mathbf{w}_t}^f(\tilde{x}, x_i)}_{\text{eNTK of backbone}} \underbrace{(\nabla_z q_t(x_i))^{\mathsf{T}}}_{\text{direction of head}} \underbrace{(e_{y_i} - q_t(x_i))}_{\text{"energy"}} + \mathcal{O}(\eta^2)$$

- If initial "energy", e.g. $\mathbb{E}_{x_i, y_i} \left\| e_{y_i} - p_0(x_i) \right\|$, is small, features don't change much

# How much do we change our features?

- Same kind of decomposition with backbone features $z = f(x)$, head $q = \mathrm{softmax}(g(z))$:

$$z_{t+1}(\tilde{x}) - z_t(\tilde{x}) = \frac{\eta}{N} \sum_{i=1}^{N} \underbrace{\mathrm{eNTK}_{\mathbf{w}_t}^f(\tilde{x}, x_i)}_{\text{eNTK of backbone}} \underbrace{(\nabla_z q_t(x_i))^\mathsf{T}}_{\text{direction of head}} \underbrace{(e_{y_i} - q_t(x_i))}_{\text{``energy''}} + \mathcal{O}(\eta^2)$$

- If initial "energy", e.g. $\mathbb{E}_{x_i, y_i} \left\| e_{y_i} - p_0(x_i) \right\|$, is small, features don't change much

- If we didn't do any head probing, "direction" is very random, especially if $g$ is rich

# How much do we change our features?

- Same kind of decomposition with backbone features $z = f(x)$, head $q = \mathrm{softmax}(g(z))$:

$$z_{t+1}(\tilde{x}) - z_t(\tilde{x}) = \frac{\eta}{N} \sum_{i=1}^{N} \underbrace{\mathrm{eNTK}^f_{\mathbf{w}_t}(\tilde{x}, x_i)}_{\text{eNTK of backbone}} \underbrace{(\nabla_z q_t(x_i))^{\mathsf{T}}}_{\text{direction of head}} \underbrace{(e_{y_i} - q_t(x_i))}_{\text{"energy"}} + \mathcal{O}(\eta^2)$$

- If initial "energy", e.g. $\mathbb{E}_{x_i, y_i} \| e_{y_i} - p_0(x_i) \|$, is small, features don't change much

- If we didn't do any head probing, "direction" is very random, especially if $g$ is rich

- Specializing to simple linear-linear model, can get insights about trends in $z$

# How much do we change our features?

- Same kind of decomposition with backbone features $z = f(x)$, head $q = \mathrm{softmax}(g(z))$:

$$z_{t+1}(\tilde{x}) - z_t(\tilde{x}) = \frac{\eta}{N} \sum_{i=1}^{N} \underbrace{\mathrm{eNTK}_{\mathbf{w}_t}^f(\tilde{x}, x_i)}_{\text{eNTK of backbone}} \underbrace{(\nabla_z q_t(x_i))^\mathsf{T}}_{\text{direction of head}} \underbrace{(e_{y_i} - q_t(x_i))}_{\text{``energy''}} + \mathcal{O}(\eta^2)$$

- If initial "energy", e.g. $\mathbb{E}_{x_i, y_i} \left\| e_{y_i} - p_0(x_i) \right\|$, is small, features don't change much

- If we didn't do any head probing, "direction" is very random, especially if $g$ is rich

- Specializing to simple linear-linear model, can get insights about trends in $z$

- Recommendations from paper:
  - Early stop during head probing (ideally, try multiple lengths for downstream task)
  - Label smoothing can help; so can more complex heads, but be careful

# How good will our fine-tuned features be? [Wei/Hu/Steinhardt 2022]

# What can we learn from empirical NTKs?

In this talk:

- As a theoretical-ish tool for local understanding:
    - Fine-grained explanation for early stopping in knowledge distillation
    - How you should fine-tune models

- As a practical tool for approximating "lookahead" in active learning

- Plus: efficiently approximating **eNTK**s for large output dimensions $k$, with guarantees

# Pool-based active learning

# Pool-based active learning

# Pool-based active learning

# Pool-based active learning

# Pool-based active learning

# Approximate retraining with local linearization

- Given $f_{\mathcal{L}}$ trained on labeled data $\mathcal{L}$, approximate $f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ with local linearization

# Approximate retraining with local linearization

- Given $f_{\mathcal{L}}$ trained on labeled data $\mathcal{L}$, approximate $f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ with local linearization

$$f_{\mathcal{L} \cup \{(x_i, y_i)\}}(\tilde{x}) \approx f_{\mathcal{L}}(\tilde{x}) + \text{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\tilde{x}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right) \text{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)^{-1} \left(\begin{bmatrix} \mathbf{y}_{\mathcal{L}} \\ y_i \end{bmatrix} - f_{\mathcal{L}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)\right)$$

# Approximate retraining with local linearization

- Given $f_{\mathcal{L}}$ trained on labeled data $\mathcal{L}$, approximate $f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ with local linearization

$$f_{\mathcal{L} \cup \{(x_i, y_i)\}}(\tilde{x}) \approx f_{\mathcal{L}}(\tilde{x}) + \text{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\tilde{x}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right) \text{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)^{-1}\left(\begin{bmatrix} \mathbf{y}_{\mathcal{L}} \\ y_i \end{bmatrix} - f_{\mathcal{L}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)\right)$$

  - Rank-one updates for efficient computation: schema

# Approximate retraining with local linearization

- Given $f_{\mathcal{L}}$ trained on labeled data $\mathcal{L}$, approximate $f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ with local linearization

$$f_{\mathcal{L} \cup \{(x_i, y_i)\}}(\tilde{x}) \approx f_{\mathcal{L}}(\tilde{x}) + \mathrm{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\tilde{x}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right) \mathrm{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)^{-1}\left(\begin{bmatrix} \mathbf{y}_{\mathcal{L}} \\ y_i \end{bmatrix} - f_{\mathcal{L}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)\right)$$

  - Rank-one updates for efficient computation: schema $\square + \boxed{\phantom{xx}}\square \times \boxed{\phantom{x}}^{-1}\left(\begin{bmatrix}\phantom{x}\end{bmatrix} - \begin{bmatrix}\phantom{x}\end{bmatrix}\right)$

- We prove this is exact for infinitely wide networks
  - $f_0 \to f_{\mathcal{L}} \to f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ agrees with direct $f_0 \to f_{\mathcal{L} \cup \{(x_i, y_i)\}}$
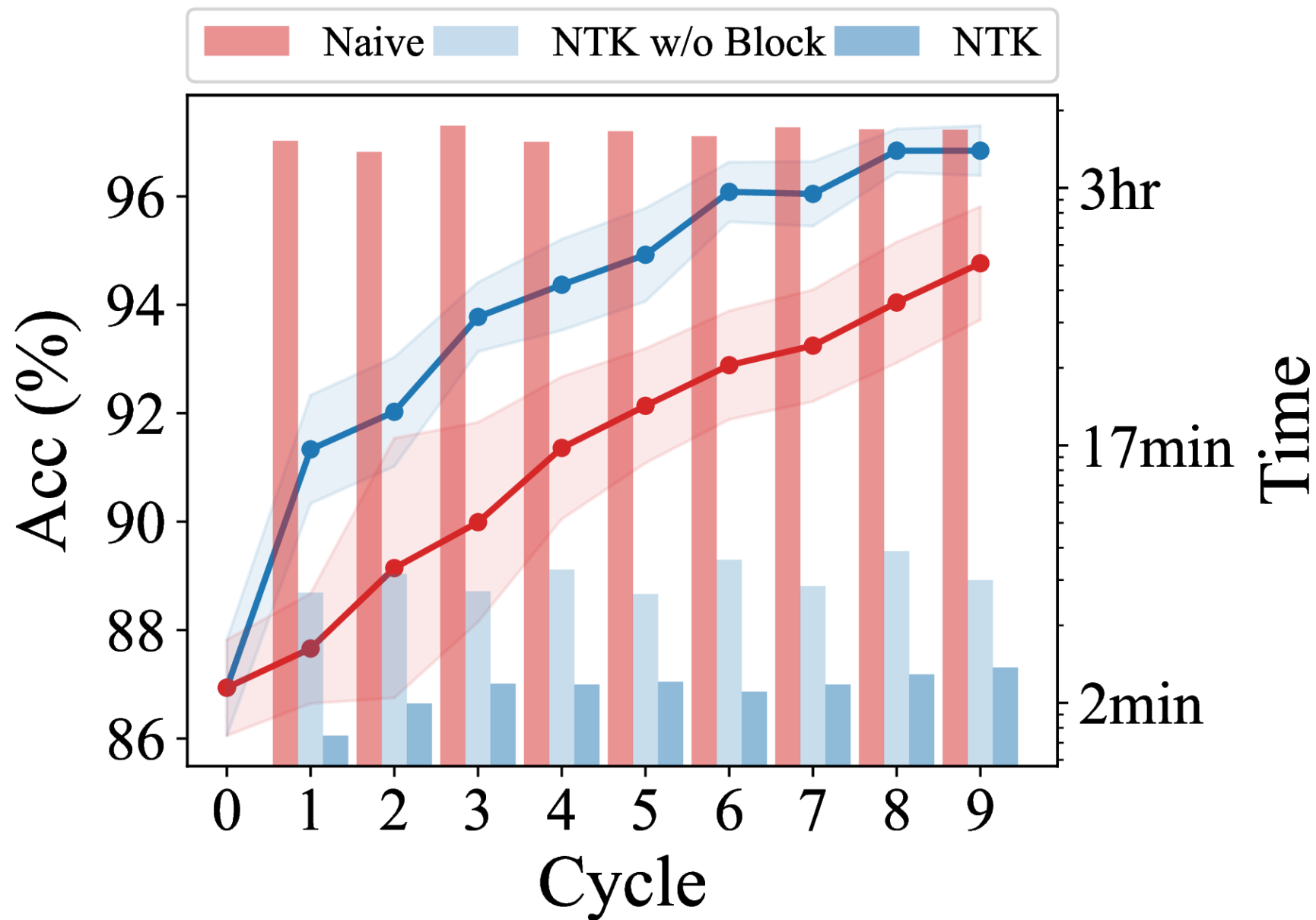
# Approximate retraining with local linearization

- Given $f_{\mathcal{L}}$ trained on labeled data $\mathcal{L}$, approximate $f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ with local linearization

$$f_{\mathcal{L} \cup \{(x_i, y_i)\}}(\tilde{x}) \approx f_{\mathcal{L}}(\tilde{x}) + \mathrm{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\tilde{x}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right) \mathrm{eNTK}_{\mathbf{w}_{\mathcal{L}}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}, \begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)^{-1} \left(\begin{bmatrix} \mathbf{y}_{\mathcal{L}} \\ y_i \end{bmatrix} - f_{\mathcal{L}}\left(\begin{bmatrix} \mathbf{X}_{\mathcal{L}} \\ x_i \end{bmatrix}\right)\right)$$
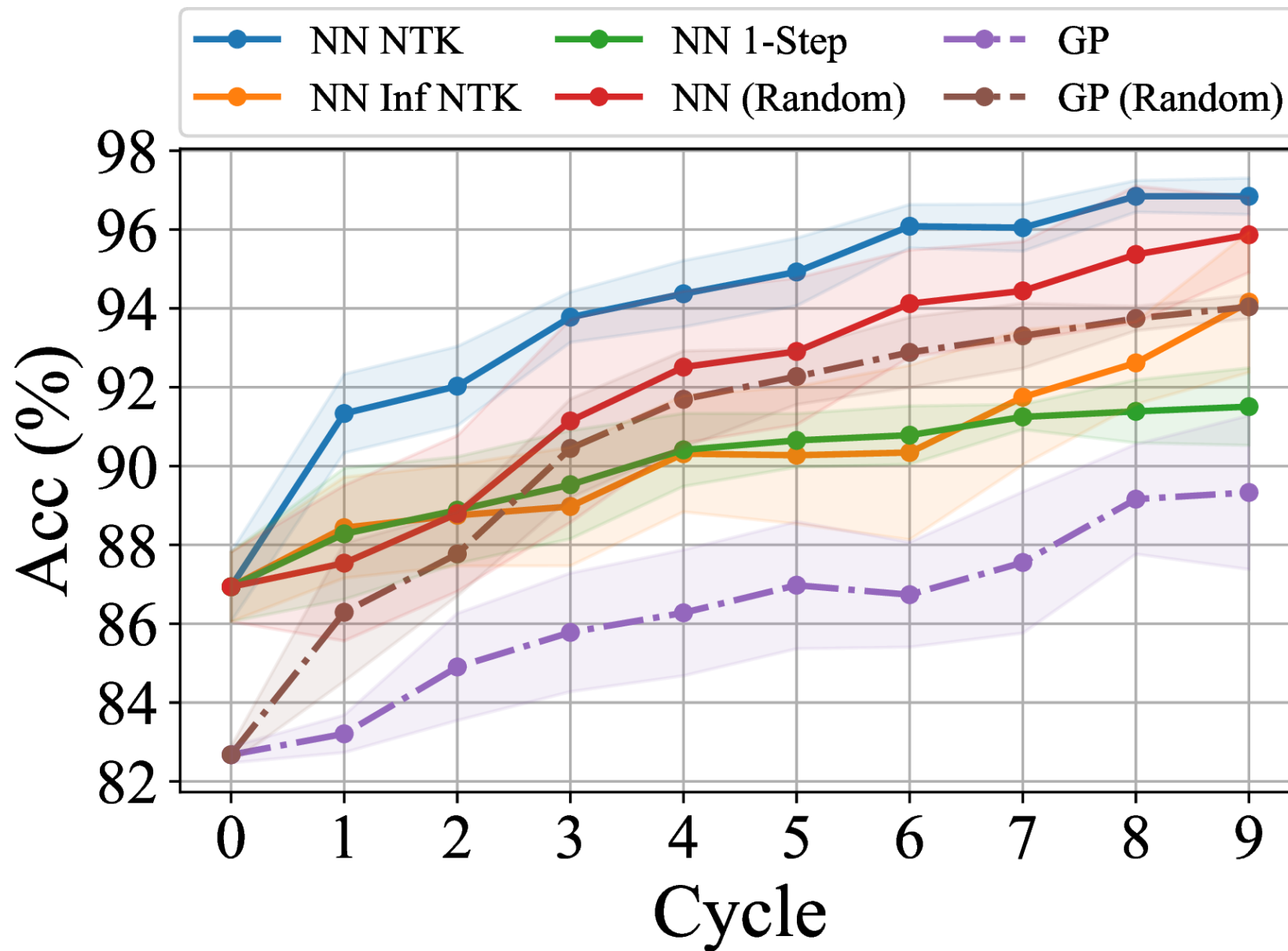
  - Rank-one updates for efficient computation: schema $\square + \boxed{\phantom{xx}}\square \times \boxed{\phantom{x}}^{-1}\left(\Big| - \Big|\right)$

- We prove this is exact for infinitely wide networks
  - $f_0 \to f_{\mathcal{L}} \to f_{\mathcal{L} \cup \{(x_i, y_i)\}}$ agrees with direct $f_0 \to f_{\mathcal{L} \cup \{(x_i, y_i)\}}$

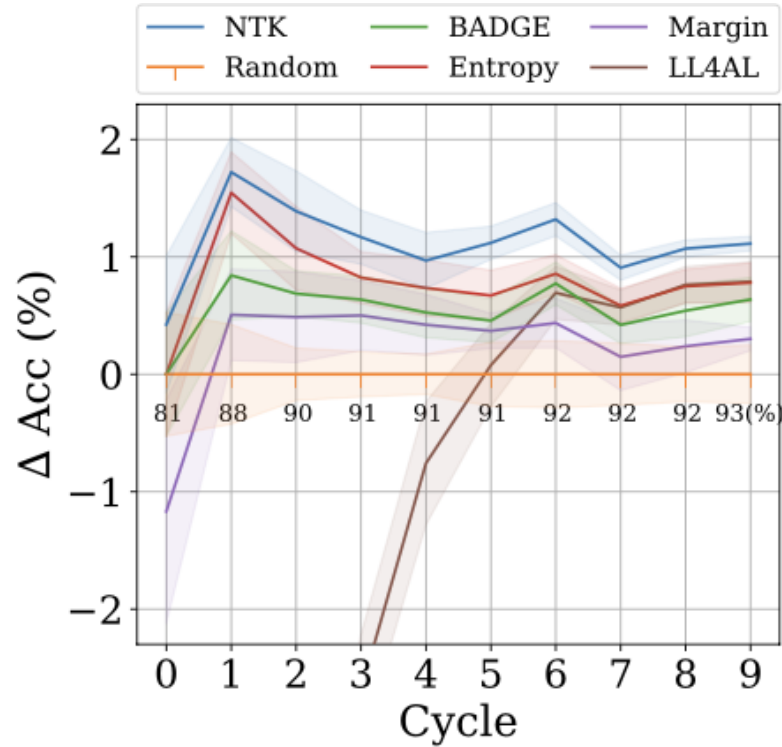- Local approximation with eNTK "should" work much more broadly than "NTK regime"
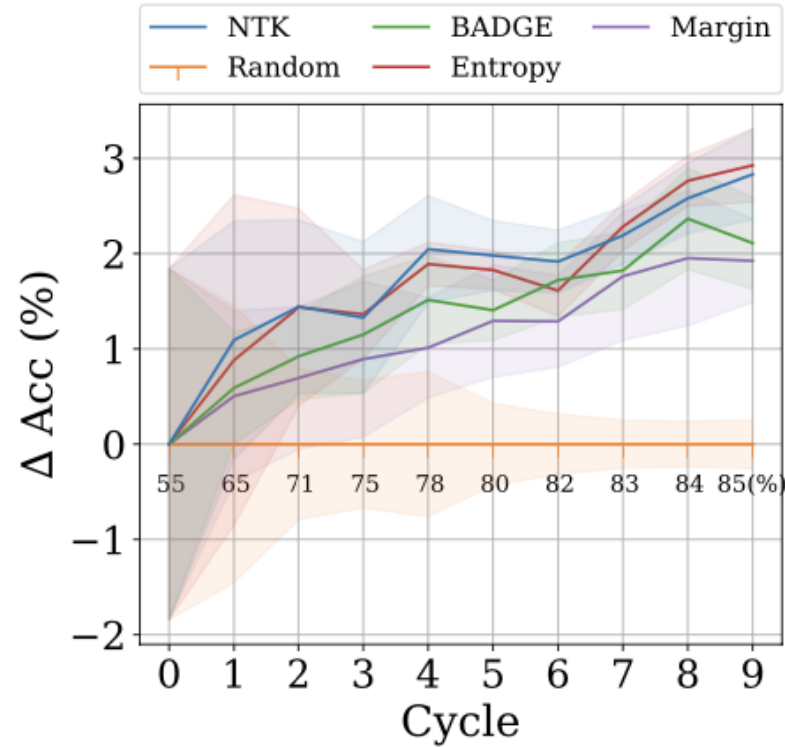
# Much faster than SGD

Much more effective than infinite NTK and one-step SGD

Legend: NN NTK, NN Inf NTK, NN 1-Step, NN (Random), GP, GP (Random)
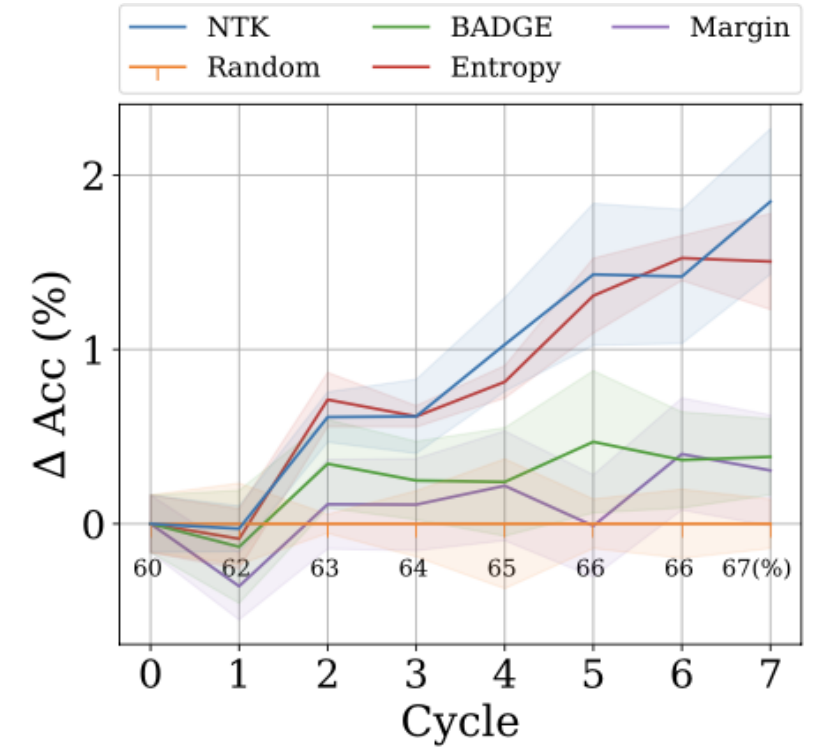
Y-axis: Acc (%)
X-axis: Cycle

# Matches/beats state of the art



(a) SVHN: 1-layer WideResNet  (b) CIFAR10: 2-layer WideResNet  (c) CIFAR100: ResNet18

Figure 2: Comparison of the-state-of-the-art active learning methods on various benchmark datasets. Vertical axis shows difference from random acquisition, whose accuracy is shown in text.

Downside: usually more computationally expensive (especially memory)

# Enables new interaction modes

# What can we learn from empirical NTKs?

In this talk:

- As a theoretical-ish tool for local understanding:
  - Fine-grained explanation for early stopping in knowledge distillation
  - How you should fine-tune models
- As a practical tool for approximating "lookahead" in active learning
- Plus: efficiently approximating **eNTK**s for large output dimensions $k$, with guarantees

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results…

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results...

- With $k$ classes, $\mathrm{eNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{kN \times kN}$ – potentially very big

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results...

- With $k$ classes, $\mathbf{eNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{kN \times kN}$ – potentially very big

- But actually, we know that $\mathbb{E}_{\mathbf{w}}\, \mathbf{eNTK_w}(x_1, x_2)$ is diagonal for most architectures
  - Let $\mathbf{pNTK_w}(x_1, x_2) = \underbrace{[\nabla_{\mathbf{w}} f_1(x_1)]}_{1 \times p}\underbrace{[\nabla_{\mathbf{w}} f_1(x_2))^\top}_{p \times 1}.$

$$\mathbb{E}_{\mathbf{w}}\, \mathbf{eNTK_w}(x_1, x_2) = \mathbb{E}_w[\mathbf{pNTK_w}(x_1, x_2)]I_k. \qquad \mathbf{pNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N} \text{ (no } k\text{!)}$$

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results...

- With $k$ classes, $\mathrm{eNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{kN \times kN}$ – potentially very big

- But actually, we know that $\mathbb{E}_{\mathbf{w}} \, \mathrm{eNTK}_{\mathbf{w}}(x_1, x_2)$ is diagonal for most architectures

  - Let $\mathrm{pNTK}_{\mathbf{w}}(x_1, x_2) = \underbrace{[\nabla_{\mathbf{w}} f_1(x_1)]}_{1 \times p} \underbrace{[\nabla_{\mathbf{w}} f_1(x_2))^{\mathsf{T}}}_{p \times 1}.$

$$\mathbb{E}_{\mathbf{w}} \, \mathrm{eNTK}_{\mathbf{w}}(x_1, x_2) = \mathbb{E}_{w}[\mathrm{pNTK}_{\mathbf{w}}(x_1, x_2)] I_k. \qquad \mathrm{pNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N} \text{ (no } k!)$$

  - Can also use "sum of logits" $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} f_j$ instead of just "first logit" $f_1$

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results...

- With $k$ classes, $\mathbf{eNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{kN \times kN}$ – potentially very big

- But actually, we know that $\mathbb{E}_{\mathbf{w}}\, \mathbf{eNTK}_{\mathbf{w}}(x_1, x_2)$ is diagonal for most architectures
  - Let $\mathbf{pNTK}_{\mathbf{w}}(x_1, x_2) = \underbrace{[\nabla_{\mathbf{w}} f_1(x_1)]}_{1 \times p}\underbrace{[\nabla_{\mathbf{w}} f_1(x_2))^{\mathsf{T}}}_{p \times 1}.$

    $$\mathbb{E}_{\mathbf{w}}\, \mathbf{eNTK}_{\mathbf{w}}(x_1, x_2) = \mathbb{E}_w[\mathbf{pNTK}_{\mathbf{w}}(x_1, x_2)]I_k. \qquad \mathbf{pNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N} \text{ (no } k!)$$

  - Can also use "sum of logits" $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} f_j$ instead of just "first logit" $f_1$

- Lots of work (including above) has used $\mathbf{pNTK}$ instead of $\mathbf{eNTK}$
  - Often without saying anything; sometimes doesn't seem like they know they're doing it

# Approximating empirical NTKs

- I hid something from you on active learning (and Wei/Hu/Steinhardt fine-tuning) results...

- With $k$ classes, $\mathbf{eNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{kN \times kN}$ – potentially very big

- But actually, we know that $\mathbb{E}_{\mathbf{w}} \, \mathbf{eNTK}_{\mathbf{w}}(x_1, x_2)$ is diagonal for most architectures
  - Let $\mathbf{pNTK}_{\mathbf{w}}(x_1, x_2) = \underbrace{[\nabla_{\mathbf{w}} f_1(x_1)]}_{1 \times p} \underbrace{[\nabla_{\mathbf{w}} f_1(x_2))^{\mathsf{T}}}_{p \times 1}.$

  $$\mathbb{E}_{\mathbf{w}} \, \mathbf{eNTK}_{\mathbf{w}}(x_1, x_2) = \mathbb{E}_w [\mathbf{pNTK}_{\mathbf{w}}(x_1, x_2)] I_k. \qquad \mathbf{pNTK}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N} \text{ (no } k!)$$

  - Can also use "sum of logits" $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} f_j$ instead of just "first logit" $f_1$

- Lots of work (including above) has used $\mathbf{pNTK}$ instead of $\mathbf{eNTK}$
  - Often without saying anything; sometimes doesn't seem like they know they're doing it

- Can we justify this more rigorously?

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} v_j$ have same distribution

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} v_j$ have same distribution

$$\mathrm{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^\mathsf{T} \, \mathrm{eNTK}^{\phi}_{\mathbf{w} \backslash V}(x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^\mathsf{T}\phi(x_2)$$

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} v_j$ have same distribution

$$\text{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^{\mathsf{T}} \, \text{eNTK}^{\phi}_{\mathbf{w} \setminus V}(x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^{\mathsf{T}}\phi(x_2)$$

$$\text{pNTK}_{\mathbf{w}}(x_1, x_2) = v_1^{\mathsf{T}} \, \text{eNTK}^{\phi}_{\mathbf{w} \setminus V}(x_1, x_2) \, v_1 + \phi(x_1)^{\mathsf{T}}\phi(x_2)$$

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^k v_j$ have same distribution

$$\text{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^\top \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi(x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^\top \phi(x_2)$$

$$\text{pNTK}_{\mathbf{w}}(x_1, x_2) = v_1^\top \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi(x_1, x_2) \, v_1 + \phi(x_1)^\top \phi(x_2)$$

- We want to bound difference $\text{eNTK}(x_1, x_2) - \text{pNTK}(x_1, x_2)I_k$

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} v_j$ have same distribution

$$\mathrm{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^\mathsf{T} \, \mathrm{eNTK}_{\mathbf{w} \backslash V}^{\phi}(x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^\mathsf{T}\phi(x_2)$$

$$\mathrm{pNTK}_{\mathbf{w}}(x_1, x_2) = v_1^\mathsf{T} \, \mathrm{eNTK}_{\mathbf{w} \backslash V}^{\phi}(x_1, x_2) \, v_1 + \phi(x_1)^\mathsf{T}\phi(x_2)$$

- We want to bound difference $\mathrm{eNTK}(x_1, x_2) - \mathrm{pNTK}(x_1, x_2)I_k$
  - Want $v_1^\mathsf{T} A v_1$ and $v_j^\mathsf{T} A v_j$ to be close, and $v_j^\mathsf{T} A v_{j'}$ small, for random $v$ and fixed $A$

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^k v_j$ have same distribution

$$\text{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^\mathsf{T} \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi(x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^\mathsf{T}\phi(x_2)$$

$$\text{pNTK}_{\mathbf{w}}(x_1, x_2) = \textcolor{red}{v_1^\mathsf{T}} \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi(x_1, x_2) \, \textcolor{red}{v_1} + \phi(x_1)^\mathsf{T}\phi(x_2)$$

- We want to bound difference $\text{eNTK}(x_1, x_2) - \text{pNTK}(x_1, x_2)I_k$
  - Want $v_1^\mathsf{T} A v_1$ and $v_j^\mathsf{T} A v_j$ to be close, and $v_j^\mathsf{T} A v_{j'}$ small, for random $v$ and fixed $A$

  - Using Hanson-Wright: $\dfrac{\|\text{eNTK} - \text{pNTK}\, I\|_F}{\|\text{eNTK}\|_F} \leq \dfrac{\|\text{eNTK}^\phi\|_F + 4\sqrt{h}}{\text{Tr}(\text{eNTK}^\phi)} k \log \dfrac{2k^2}{\delta}$

# pNTK motivation

- Say $f(x) = V\phi(x)$, $\phi(x) \in \mathbb{R}^h$, and $V \in \mathbb{R}^{k \times h}$ has rows $v_j \in \mathbb{R}^h$ with iid entries

- If $v_{j,i} \sim \mathcal{N}(0, \sigma^2)$, then $v_1$ and $\frac{1}{\sqrt{k}} \sum_{j=1}^{k} v_j$ have same distribution

$$\text{eNTK}_{\mathbf{w}}(x_1, x_2)_{jj'} = v_j^\mathsf{T} \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi (x_1, x_2) \, v_{j'} + \mathbb{I}(j = j')\phi(x_1)^\mathsf{T} \phi(x_2)$$

$$\text{pNTK}_{\mathbf{w}}(x_1, x_2) = {\color{red}v_1^\mathsf{T}} \, \text{eNTK}_{\mathbf{w} \backslash V}^\phi (x_1, x_2) \, {\color{red}v_1} + \phi(x_1)^\mathsf{T} \phi(x_2)$$

- We want to bound difference $\text{eNTK}(x_1, x_2) - \text{pNTK}(x_1, x_2)I_k$
  - Want $v_1^\mathsf{T} A v_1$ and $v_j^\mathsf{T} A v_j$ to be close, and $v_j^\mathsf{T} A v_{j'}$ small, for random $v$ and fixed $A$

  - Using Hanson-Wright: $\dfrac{\left\| \text{eNTK} - \text{pNTK} \, I \right\|_F}{\left\| \text{eNTK} \right\|_F} \leq \dfrac{\left\| \text{eNTK}^\phi \right\|_F + 4\sqrt{h}}{\text{Tr}(\text{eNTK}^\phi)} k \log \dfrac{2k^2}{\delta}$

  - Fully-connected ReLU nets at init., fan-in mode: numerator $\mathcal{O}(h\sqrt{h})$, denom $\Theta(h^2)$
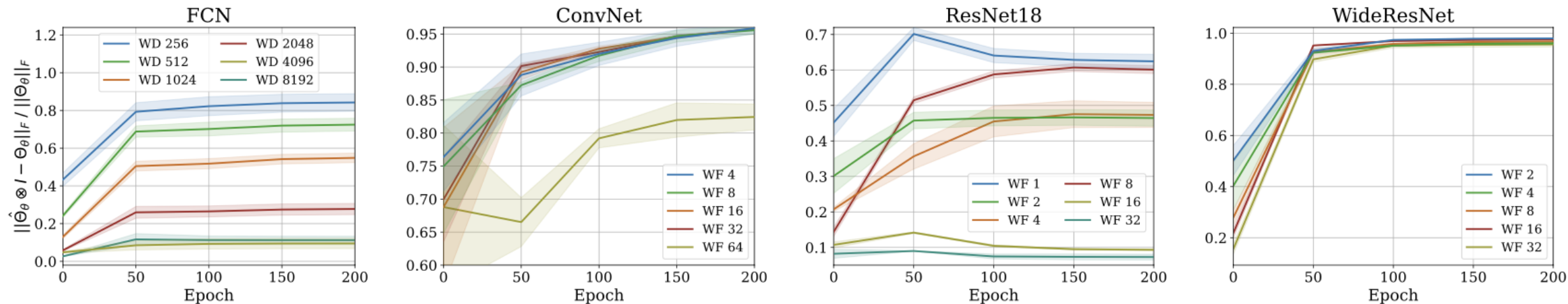
# pNTK's Frobenius error



Figure 3: Evaluating the **relative difference of Frobenius norm of** $\Theta_\theta(\mathcal{D}, \mathcal{D})$ **and** $\hat{\Theta}_\theta(\mathcal{D}, \mathcal{D}) \otimes I_O$ at initialization and throughout training, based on $\mathcal{D}$ being 1000 random points from CIFAR-10. Wider nets have more similar $\|\Theta_\theta\|_F$ and $\|\hat{\Theta}_\theta \otimes I_O\|_F$ at initialization.
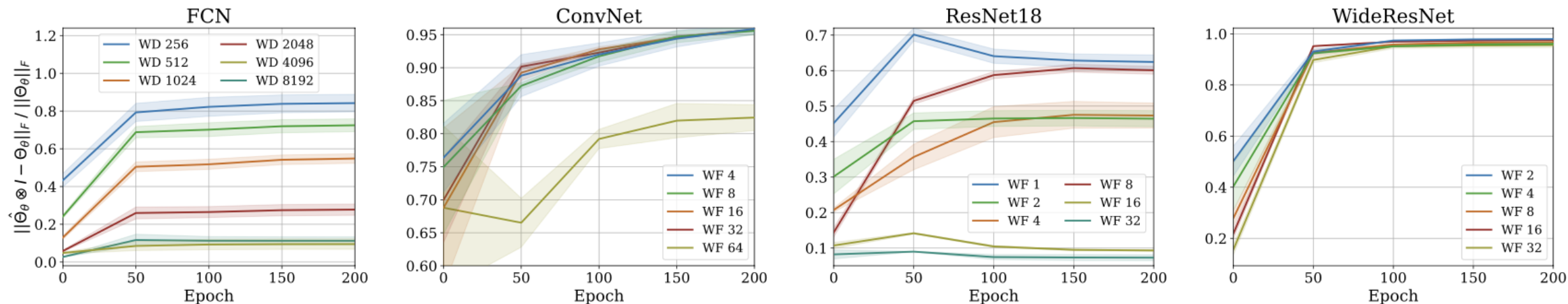
# pNTK's Frobenius error



Figure 3: Evaluating the **relative difference of Frobenius norm of** $\Theta_\theta(\mathcal{D}, \mathcal{D})$ **and** $\hat{\Theta}_\theta(\mathcal{D}, \mathcal{D}) \otimes I_O$ at initialization and throughout training, based on $\mathcal{D}$ being 1000 random points from CIFAR-10. Wider nets have more similar $\|\Theta_\theta\|_F$ and $\|\hat{\Theta}_\theta \otimes I_O\|_F$ at initialization.

Same kind of theorem / empirical results for largest eigenvalue,
and empirical results for $\lambda_{\min}$, condition number

1

# Kernel regression with pNTK

- Reshape things to handle prediction appropriately:

$$\underbrace{f_{\mathrm{eNTK}}(\tilde{x})}_{k \times 1} = \underbrace{f_0(\tilde{x})}_{k \times 1} + \underbrace{\mathrm{eNTK}_{\mathbf{w}_0}(\tilde{x}, \mathbf{X})}_{k \times kN} \underbrace{\mathrm{eNTK}_{\mathbf{w}_0}(\mathbf{X}, \mathbf{X})^{-1}}_{kN \times kN} \underbrace{(\mathbf{y} - f_0(\mathbf{X}))}_{kN \times 1}$$

$$\underbrace{f_{\mathrm{pNTK}}(\tilde{x})}_{k \times 1} = \underbrace{f_0(\tilde{x})}_{k \times 1} + \left( \underbrace{\mathrm{pNTK}_{\mathbf{w}_0}(\tilde{x}, \mathbf{X})}_{1 \times N} \underbrace{\mathrm{pNTK}_{\mathbf{w}_0}(\mathbf{X}, \mathbf{X})^{-1}}_{N \times N} \underbrace{(\mathbf{y} - f_0(\mathbf{X}))}_{N \times k} \right)^{\mathsf{T}}$$

- We have $\|f_{\mathrm{eNTK}}(\tilde{x}) - f_{\mathrm{pNTK}}(\tilde{x})\| = \mathcal{O}(\frac{1}{\sqrt{h}})$ again

  - If we add regularization, need to "scale" $\lambda$ between the two
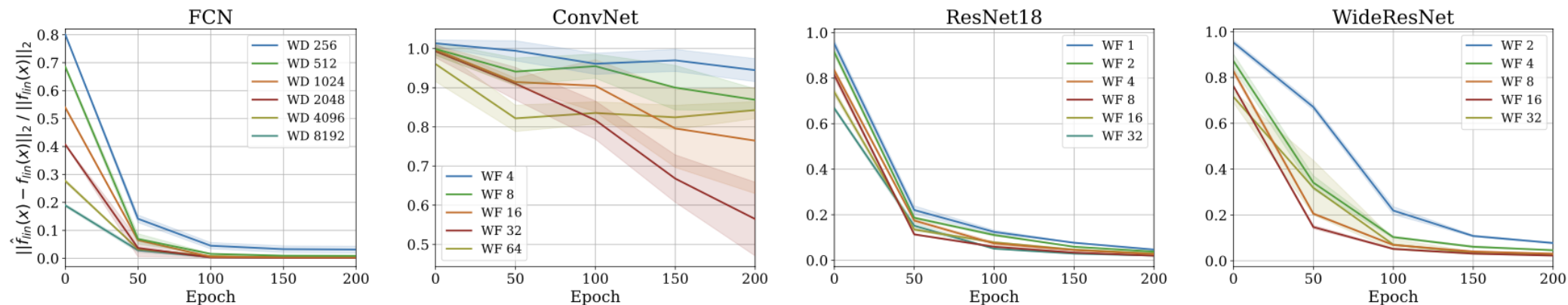
# Kernel regression with pNTK



Figure 7: The **relative difference of kernel regression outputs**, (4) and (5), when training on $|\mathcal{D}| = 1000$ random CIFAR-10 points and testing on $|\mathcal{X}| = 500$. For wider NNs, the relative difference in $\hat{f}^{lin}(\mathcal{X})$ and $f^{lin}(\mathcal{X})$ decreases at initialization. Surprisingly, the difference between these two continues to quickly vanish while training the network.
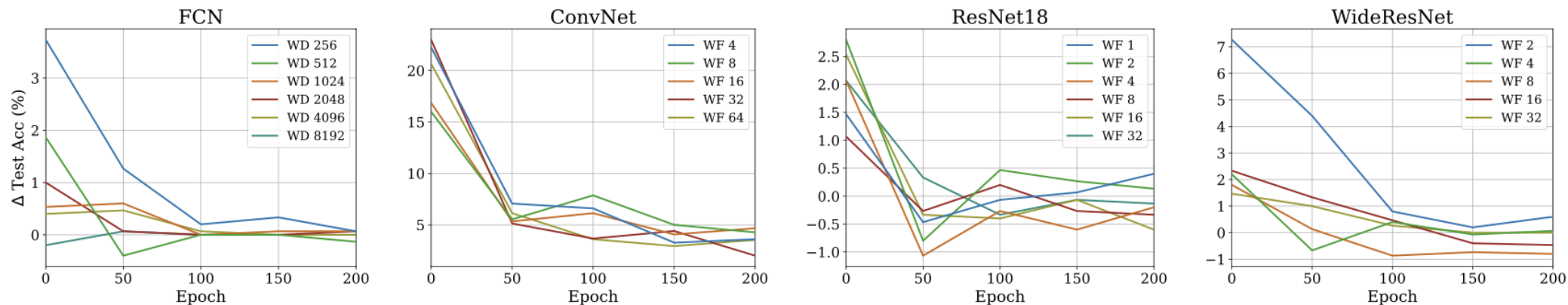
Figure 8: **Using pNTK in kernel regression** (as in Figure 7) **almost always achieves a higher test accuracy than using eNTK.** Wider NNs and trained nets have more similar prediction accuracies of $\hat{f}^{lin}$ and $f^{lin}$ at initialization. Again, the difference between these two continues to vanish throughout the training process using SGD.
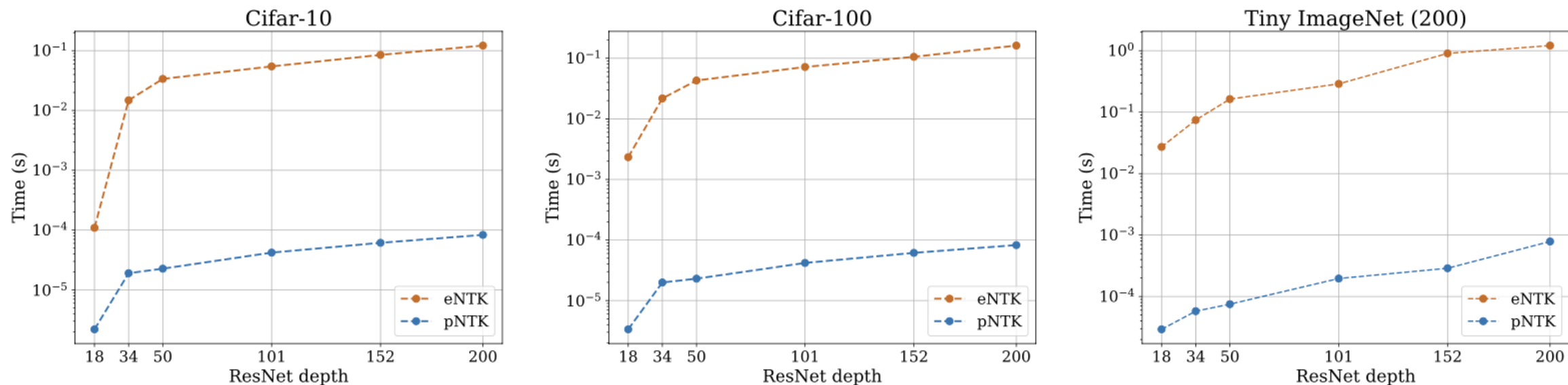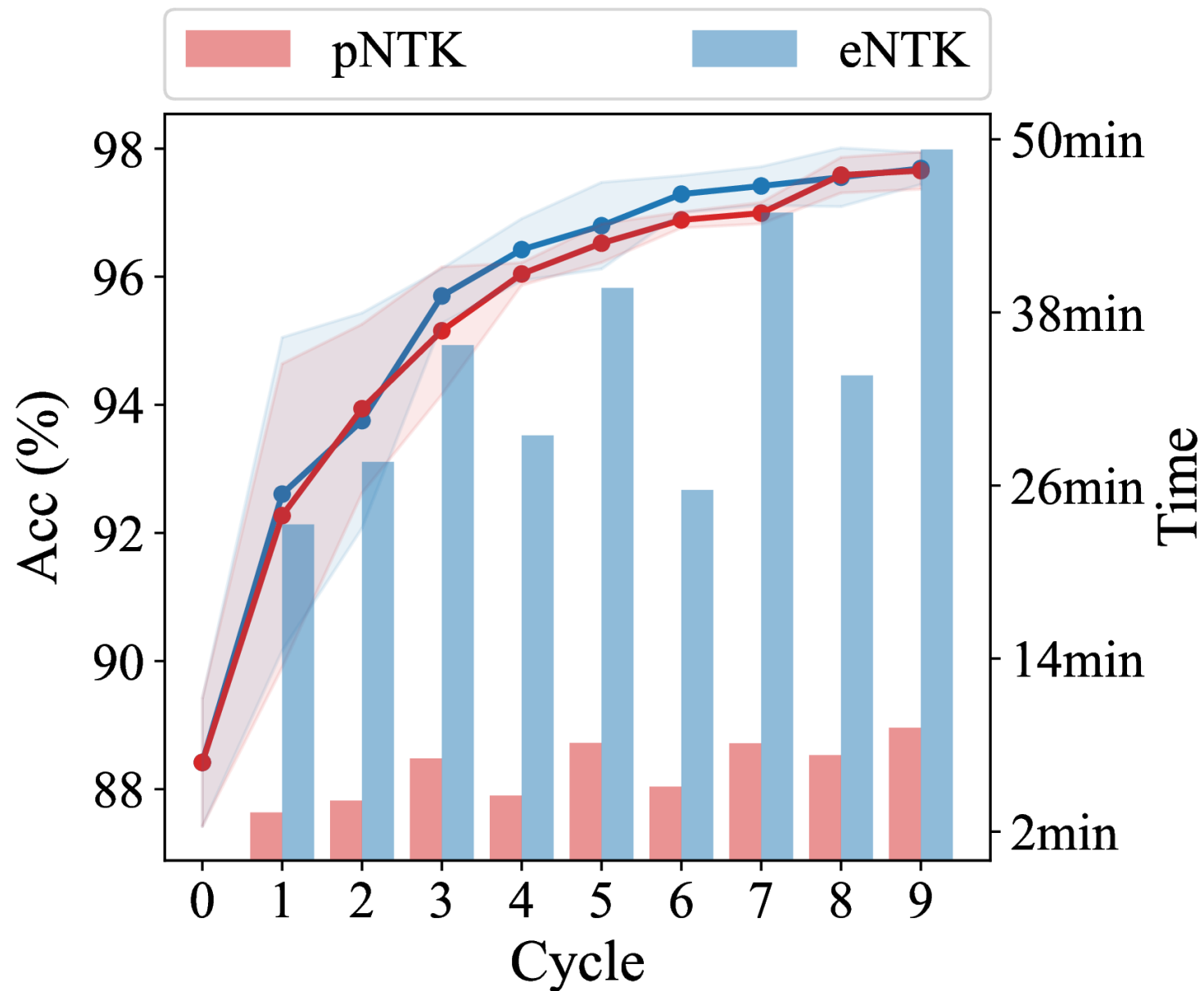
# pNTK speed-up



Figure 1: **Wall-clock time to evaluate the eNTK and pNTK** for one pair of inputs, across datasets and ResNet depths.

pNTK speed-up on active learning task

# pNTK for full CIFAR-10 regression

- $\mathrm{eNTK}(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 1.8 terabytes of memory

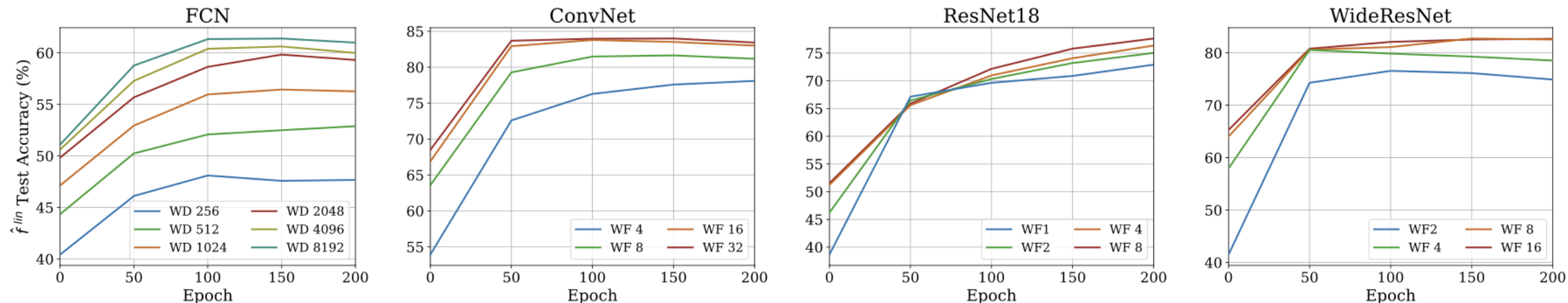- $\mathrm{pNTK}(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 18 gigabytes of memory



Figure 9: Evaluating the **test accuracy of kernel regression predictions using pNTK as in (5) on the full CIFAR-10 dataset**. As the NN's width grows, the test accuracy of $\hat{f}^{lin}$ also improves, but eventually saturates with the growing width. Using trained weights in computation of pNTK results in improved test accuracy of $\hat{f}^{lin}$.

# pNTK for full CIFAR-10 regression

- $\mathrm{eNTK}(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 1.8 terabytes of memory

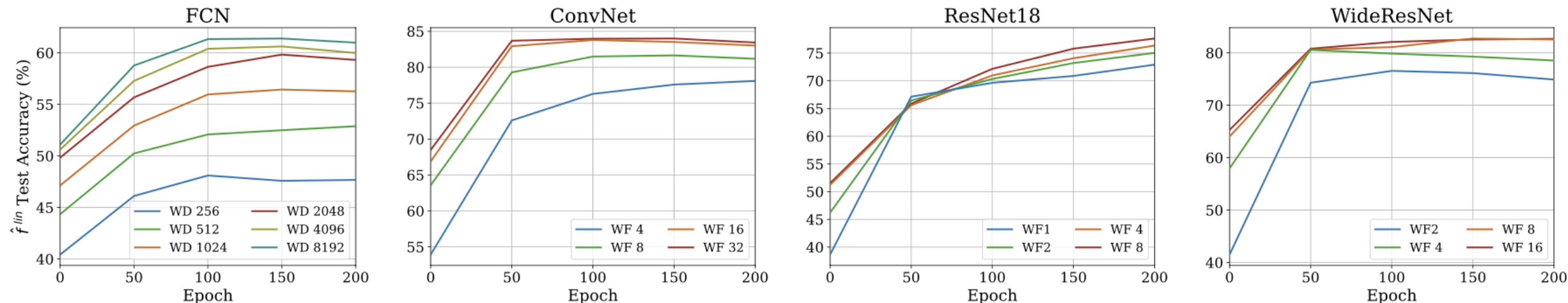- $\mathrm{pNTK}(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 18 gigabytes of memory



Figure 9: Evaluating the **test accuracy of kernel regression predictions using pNTK as in (5) on the full CIFAR-10 dataset**. As the NN's width grows, the test accuracy of $\hat{f}^{lin}$ also improves, but eventually saturates with the growing width. Using trained weights in computation of pNTK results in improved test accuracy of $\hat{f}^{lin}$.

- Worse than infinite NTK for FCN/ConvNet (where they can be computed, if you try hard)

# pNTK for full CIFAR-10 regression

- $eNTK(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 1.8 terabytes of memory

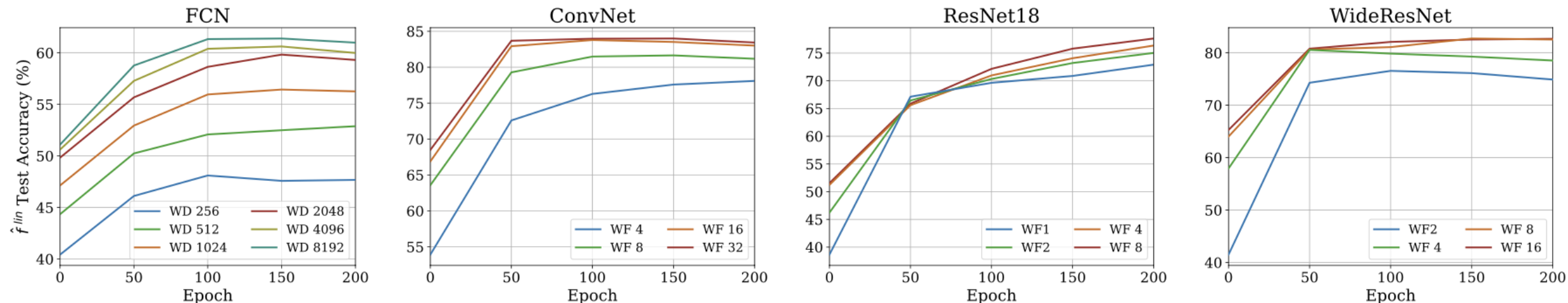- $pNTK(\mathbf{X}, \mathbf{X})$ on CIFAR-10: 18 gigabytes of memory



Figure 9: Evaluating the **test accuracy of kernel regression predictions using pNTK as in (5) on the full CIFAR-10 dataset**. As the NN's width grows, the test accuracy of $\hat{f}^{lin}$ also improves, but eventually saturates with the growing width. Using trained weights in computation of pNTK results in improved test accuracy of $\hat{f}^{lin}$.

- Worse than infinite NTK for FCN/ConvNet (where they can be computed, if you try hard)

- Way worse than SGD

# Recap

eNTK is a good tool for intuitive understanding of the learning process

| |
|---|
| Ren, Guo, S.    Better Supervisory Signals by Observing Learning Paths |

| |
|---|
| Ren, Guo, Bae, S.    How to prepare your task head for finetuning |

## eNTK is practically very effective at "lookahead" for active learning

| |
|---|
| Mohamadi*, Bae*, S.    Making Look-Ahead Active Learning Strategies Feasible with Neural Tangent Kernels |

## You should probably use pNTK instead of eNTK for high-dim output problems:

| |
|---|
| Mohamadi, Bae, S.    A Fast, Well-Founded Approximation to the Empirical Neural Tangent Kernel |