

# Infinite Width, NTK, and Feature Learning

Student Meeting 1

---

Shreyas Kalvankar

September 17, 2025

TU Delft

# Agenda

## Goal of this talk

- Introduce the infinite-width lens, NTK, and the lazy training regime via a simple one-layer example.
- Situate these ideas within related work: baseline results at infinite width and viewpoints that go beyond linearization.

## What I would like from you

- Conceptual clarifications and critiques of how I'm connecting the papers.
- Suggestions for key references I might be missing.
- Pointers to alternative frameworks worth comparing (mean-field, functional-analytic, optimization bias).

# Why theory?

- Larger models keep improving, but we don't fully know *why*.

# Why theory?

- Larger models keep improving, but we don't fully know *why*.
- **A solvable starting point:** infinite width gives a clean baseline; modern models are large-but-finite, deep, and do feature learning.

# Why theory?

- Larger models keep improving, but we don't fully know *why*.
- **A solvable starting point:** infinite width gives a clean baseline; modern models are large-but-finite, deep, and do feature learning.
- **What can theory help in:** explanations of convergence/generalization, signals for when features move, and guidance on design choices (init, learning rate, normalization).

- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.

- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.
- **NTK:** first-order (linear) view of training; clear convergence intuition via a fixed kernel.

- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.
- **NTK:** first-order (linear) view of training; clear convergence intuition via a fixed kernel.
- **Beyond the linear NTK picture:** when does it stop being accurate?
  - Does the *kernel* change during training?



- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.
- **NTK:** first-order (linear) view of training; clear convergence intuition via a fixed kernel.
- **Beyond the linear NTK picture:** when does it stop being accurate?
  - Does the *kernel* change during training?
  - Do the model's internal features move?

- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.
- **NTK:** first-order (linear) view of training; clear convergence intuition via a fixed kernel.
- **Beyond the linear NTK picture:** when does it stop being accurate?
  - Does the *kernel* change during training?
  - Do the model's internal features move?
  - Do training curves deviate from the NTK baseline prediction?

- **Infinite width:** clean, analyzable baseline; randomness averages to a deterministic kernel.
- **NTK:** first-order (linear) view of training; clear convergence intuition via a fixed kernel.
- **Beyond the linear NTK picture:** when does it stop being accurate?
  - Does the *kernel* change during training?
  - Do the model's internal features move?
  - Do training curves deviate from the NTK baseline prediction?
  - Possible causes: higher-order effects, useful parameter re-organization that push the model out of the lazy regime.

**Supervised learning:** data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ .

**Supervised learning:** data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ .

Neural network  $f : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}$  with parameters  $\boldsymbol{\theta} \in \mathbb{R}^p$ .

**Supervised learning:** data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ .

Neural network  $f : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}$  with parameters  $\boldsymbol{\theta} \in \mathbb{R}^p$ .

Trained by (continuous-time) gradient flow on squared loss:

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

**Supervised learning:** data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ .

Neural network  $f : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}$  with parameters  $\boldsymbol{\theta} \in \mathbb{R}^p$ .

Trained by (continuous-time) gradient flow on squared loss:

$$L(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i; \boldsymbol{\theta}) - y_i)^2, \quad \frac{d}{dt} \boldsymbol{\theta}_t = -\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t).$$

# Linearization at Initialization

First-order Taylor around  $\theta_0$ :

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \underbrace{\nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}_0)^\top}_{:= \boldsymbol{\phi}(\mathbf{x})^\top} (\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$



# Linearization at Initialization

First-order Taylor around  $\theta_0$ :

$$f(\mathbf{x}; \theta) \approx f(\mathbf{x}; \theta_0) + \underbrace{\nabla_{\theta} f(\mathbf{x}; \theta_0)^{\top}}_{:= \phi(\mathbf{x})^{\top}} (\theta - \theta_0).$$

This induces the **Neural Tangent Kernel (NTK)** at  $\theta_0$ :

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^{\top} \phi(\mathbf{x}') = \nabla_{\theta} f(\mathbf{x}; \theta_0)^{\top} \nabla_{\theta} f(\mathbf{x}'; \theta_0).$$

*Ref:* (Jacot et al., 2018)

## Illustrative Model (One Hidden Layer, No Biases)

Consider a  $m$ -width neural network

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}),$$

## Illustrative Model (One Hidden Layer, No Biases)

Consider a  $m$ -width neural network

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad a_r \sim \mathcal{N}(0, \sigma_a^2), \quad \mathbf{w}_r \sim \mathcal{N}\left(0, \frac{\sigma_w^2}{d} I\right).$$

# Illustrative Model (One Hidden Layer, No Biases)

Consider a  $m$ -width neural network

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad a_r \sim \mathcal{N}(0, \sigma_a^2), \quad \mathbf{w}_r \sim \mathcal{N}\left(0, \frac{\sigma_w^2}{d} I\right).$$

Gradients:

$$\nabla_{a_r} f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sigma(\mathbf{w}_r^\top \mathbf{x}),$$

# Illustrative Model (One Hidden Layer, No Biases)

Consider a  $m$ -width neural network

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad a_r \sim \mathcal{N}(0, \sigma_a^2), \quad \mathbf{w}_r \sim \mathcal{N}\left(0, \frac{\sigma_w^2}{d} I\right).$$

Gradients:

$$\nabla_{a_r} f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad \nabla_{\mathbf{w}_r} f(\mathbf{x}) = \frac{1}{\sqrt{m}} a_r \sigma'(\mathbf{w}_r^\top \mathbf{x}) \mathbf{x}.$$

# Illustrative Model (One Hidden Layer, No Biases)

Consider a  $m$ -width neural network

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad a_r \sim \mathcal{N}(0, \sigma_a^2), \quad \mathbf{w}_r \sim \mathcal{N}\left(0, \frac{\sigma_w^2}{d} I\right).$$

Gradients:

$$\nabla_{a_r} f(\mathbf{x}) = \frac{1}{\sqrt{m}} \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad \nabla_{\mathbf{w}_r} f(\mathbf{x}) = \frac{1}{\sqrt{m}} a_r \sigma'(\mathbf{w}_r^\top \mathbf{x}) \mathbf{x}.$$

Finite-width NTK:

$$K_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{r=1}^m \sigma(\mathbf{w}_r^\top \mathbf{x}) \sigma(\mathbf{w}_r^\top \mathbf{x}') + \frac{1}{m} \sum_{r=1}^m a_r^2 \sigma'(\mathbf{w}_r^\top \mathbf{x}) \sigma'(\mathbf{w}_r^\top \mathbf{x}') \mathbf{x}^\top \mathbf{x}'.$$

## Infinite-Width NTK

The two sums are empirical averages of i.i.d. terms. Since  $a_r$  and  $w_r$  are independent with finite moments, the (strong) law of large numbers gives, almost surely, as width  $m \rightarrow \infty$ :

# Infinite-Width NTK

The two sums are empirical averages of i.i.d. terms. Since  $a_r$  and  $w_r$  are independent with finite moments, the (strong) law of large numbers gives, almost surely, as width  $m \rightarrow \infty$ :

$$\frac{1}{m} \sum_{r=1}^m \sigma(w_r^\top x) \sigma(w_r^\top x') \longrightarrow \mathbb{E}_w [\sigma(w^\top x) \sigma(w^\top x')],$$



# Infinite-Width NTK

The two sums are empirical averages of i.i.d. terms. Since  $a_r$  and  $w_r$  are independent with finite moments, the (strong) law of large numbers gives, almost surely, as width  $m \rightarrow \infty$ :

$$\frac{1}{m} \sum_{r=1}^m \sigma(w_r^\top x) \sigma(w_r^\top x') \longrightarrow \mathbb{E}_w [\sigma(w^\top x) \sigma(w^\top x')],$$

$$\frac{1}{m} \sum_{r=1}^m a_r^2 \sigma'(w_r^\top x) \sigma'(w_r^\top x') \longrightarrow \sigma_a^2 \mathbb{E}_w [\sigma'(w^\top x) \sigma'(w^\top x')].$$

# Infinite-Width NTK

The two sums are empirical averages of i.i.d. terms. Since  $a_r$  and  $w_r$  are independent with finite moments, the (strong) law of large numbers gives, almost surely, as width  $m \rightarrow \infty$ :

$$\frac{1}{m} \sum_{r=1}^m \sigma(w_r^\top x) \sigma(w_r^\top x') \longrightarrow \mathbb{E}_w [\sigma(w^\top x) \sigma(w^\top x')],$$

$$\frac{1}{m} \sum_{r=1}^m a_r^2 \sigma'(w_r^\top x) \sigma'(w_r^\top x') \longrightarrow \sigma_a^2 \mathbb{E}_w [\sigma'(w^\top x) \sigma'(w^\top x')].$$

Thus, in the infinite-width limit, the empirical NTK converges almost surely to a deterministic kernel

$$K_\infty(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w}} [\sigma(\mathbf{w}^\top \mathbf{x}) \sigma(\mathbf{w}^\top \mathbf{x}')] + \sigma_a^2 \mathbf{x}^\top \mathbf{x}' \mathbb{E}_{\mathbf{w}} [\sigma'(\mathbf{w}^\top \mathbf{x}) \sigma'(\mathbf{w}^\top \mathbf{x}')].$$

# Infinite-Width NTK

The two sums are empirical averages of i.i.d. terms. Since  $a_r$  and  $w_r$  are independent with finite moments, the (strong) law of large numbers gives, almost surely, as width  $m \rightarrow \infty$ :

$$\frac{1}{m} \sum_{r=1}^m \sigma(w_r^\top x) \sigma(w_r^\top x') \longrightarrow \mathbb{E}_w [\sigma(w^\top x) \sigma(w^\top x')],$$

$$\frac{1}{m} \sum_{r=1}^m a_r^2 \sigma'(w_r^\top x) \sigma'(w_r^\top x') \longrightarrow \sigma_a^2 \mathbb{E}_w [\sigma'(w^\top x) \sigma'(w^\top x')].$$

Thus, in the infinite-width limit, the empirical NTK converges almost surely to a deterministic kernel

$$K_\infty(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w}} [\sigma(\mathbf{w}^\top \mathbf{x}) \sigma(\mathbf{w}^\top \mathbf{x}')] + \sigma_a^2 \mathbf{x}^\top \mathbf{x}' \mathbb{E}_{\mathbf{w}} [\sigma'(\mathbf{w}^\top \mathbf{x}) \sigma'(\mathbf{w}^\top \mathbf{x}')].$$

Key point: in the wide limit,  $K_t \approx K_0$  remains *essentially constant* during training (lazy regime).

Refs: (Neal, 1996; Lee et al., 2019)

**Gradient flow & chain rule:** Let  $f_t(x_i) := f(x_i; \theta_t)$ . Then

$$\frac{d}{dt} f_t(x_i) = \nabla_{\theta} f(x_i; \theta_t)^{\top} \dot{\theta}_t = -\nabla_{\theta} f(x_i; \theta_t)^{\top} \nabla_{\theta} L(\theta_t).$$

## Training dynamics: mini-derivation

**Gradient flow & chain rule:** Let  $f_t(x_i) := f(x_i; \theta_t)$ . Then

$$\frac{d}{dt} f_t(x_i) = \nabla_{\theta} f(x_i; \theta_t)^{\top} \dot{\theta}_t = -\nabla_{\theta} f(x_i; \theta_t)^{\top} \nabla_{\theta} L(\theta_t).$$

**Loss gradient (squared loss):**

$$L(\theta) = \frac{1}{2} \sum_{j=1}^n (f(x_j; \theta) - y_j)^2, \quad \nabla_{\theta} L(\theta_t) = \sum_{j=1}^n (f_t(x_j) - y_j) \nabla_{\theta} f(x_j; \theta_t).$$

## Training dynamics: constant kernel

Substitute and define the (time-dependent) NTK.

$$\frac{d}{dt}f_t(x_i) = - \sum_{j=1}^n \underbrace{\nabla_{\theta} f(x_i; \theta_t)^{\top} \nabla_{\theta} f(x_j; \theta_t)}_{=: K_t(x_i, x_j)} (f_t(x_j) - y_j).$$

# Training dynamics: constant kernel

Substitute and define the (time-dependent) NTK.

$$\frac{d}{dt}f_t(x_i) = - \sum_{j=1}^n \underbrace{\nabla_{\theta} f(x_i; \theta_t)^{\top} \nabla_{\theta} f(x_j; \theta_t)}_{=: K_t(x_i, x_j)} (f_t(x_j) - y_j).$$

Vectorized form:

$$\dot{\hat{f}}_t = -K_t(\hat{f}_t - y).$$

# Training dynamics: constant kernel

Substitute and define the (time-dependent) NTK.

$$\frac{d}{dt}f_t(x_i) = - \sum_{j=1}^n \underbrace{\nabla_{\theta} f(x_i; \theta_t)^{\top} \nabla_{\theta} f(x_j; \theta_t)}_{=: K_t(x_i, x_j)} (f_t(x_j) - y_j).$$

Vectorized form:

$$\dot{\hat{f}}_t = -K_t (\hat{f}_t - y).$$

**Constant-kernel (NTK) regime.** If  $K_t \approx K_0 \equiv K$  (infinite width / lazy),

$$\dot{\hat{f}}_t = -K(\hat{f}_t - y) \quad \Rightarrow \quad \hat{f}_t = y + e^{-Kt} (\hat{f}_0 - y).$$

*Refs:* (Jacot et al., 2018; Lee et al., 2019)



**Lazy regime:** parameter drift is small,  $\|\theta_t - \theta_0\| \ll \|\theta_0\|$ .

- Features  $\phi(\mathbf{x})$  and kernel  $K$  stay (approximately) constant.
- Training reduces to kernel regression with fixed  $K$ .

**Limitation:** suppresses *feature learning* (representation change).

**Aim:** quantify *when* lazy holds/breaks and *how* to model beyond it.

# Diagnosing the transition to feature learning

## Some signals to consider

- **Kernel drift:** Does the NTK matrix  $K_t$  change during training? Compare  $K_t$  to  $K_0$  on the (same) data. Bigger change  $\Rightarrow$  more feature learning.

# Diagnosing the transition to feature learning

## Some signals to consider

- **Kernel drift:** Does the NTK matrix  $K_t$  change during training? Compare  $K_t$  to  $K_0$  on the (same) data. Bigger change  $\Rightarrow$  more feature learning.
- **Feature drift:** Do the tangent features  $\phi_t(x) = \nabla_{\theta} f(x; \theta_t)$  move on a small fixed set  $\mathcal{S}$ ? Track the average change from  $t = 0$ .

# Diagnosing the transition to feature learning

## Some signals to consider

- **Kernel drift:** Does the NTK matrix  $K_t$  change during training? Compare  $K_t$  to  $K_0$  on the (same) data. Bigger change  $\Rightarrow$  more feature learning.
- **Feature drift:** Do the tangent features  $\phi_t(x) = \nabla_{\theta} f(x; \theta_t)$  move on a small fixed set  $\mathcal{S}$ ? Track the average change from  $t = 0$ .

## How to (often) push out of lazy

- larger learning rate
- smaller width
- more depth / biases / normalization

## Beyond linearization I: quadratic / higher-order

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \phi(\mathbf{x})^\top \Delta \boldsymbol{\theta} + \frac{1}{2} \Delta \boldsymbol{\theta}^\top H_f(\mathbf{x}) \Delta \boldsymbol{\theta}, \quad \Delta \boldsymbol{\theta} = \boldsymbol{\theta} - \boldsymbol{\theta}_0.$$

- **Mechanism (Bai and Lee (2020)):** construct regimes where the linear term is suppressed so the *quadratic* term governs the dynamics; extendable to  $k > 2$  (“higher-order NTKs”).
- **Findings:** with the linear term suppressed, progress comes from *feature changes*; this adaptive regime is easy to optimize and can beat NTK in sample use on simple tasks.

See: “Beyond Linearization: On Quadratic and Higher-Order Approximation of Wide Neural Networks.”

## Beyond linearization II: adaptive / time-varying kernels

- Zhang et al. (2024) replace the fixed NTK with a *time-varying* kernel  $K_t$  that evolves during training (“kernel drift”).
- Features adapt during training and increasingly align with label-relevant directions (growing alignment).
- They provide a prototype of an over-parameterized Gaussian sequence model to analyze feature learning beyond the NTK picture.

See: “Towards a Statistical Understanding of Neural Networks: Beyond the NTK Theories.”

## Summary & discussion

- Infinite width as a clean baseline; NTK via linearization at initialization.
- Constant-kernel training dynamics:  $\dot{f}_t = -K(f_t - y)$  with solution  $f_t = y + e^{-Kt}(f_0 - y)$ .
- Lazy training  $\Rightarrow$  features (and  $K$ ) stay essentially fixed.
- How to spot leaving lazy: *kernel drift* ( $K_t \neq K_0$ ) and *feature drift* on a probe set.
- Beyond NTK in the literature:
  - *Quadratic / higher-order* near init (Bai and Lee, 2020).
  - *Adaptive / time-varying kernels* and alignment (Zhang et al., 2024).

## References

---

- Y. Bai and J. D. Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. In *International Conference on Learning Representations*, 2020.
- A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- J. Lee, L. Xiao, S. S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in Neural Information Processing Systems*, 2019.
- R. M. Neal. *Priors for infinite networks*. PhD thesis, University of Toronto, 1996.
- H. Zhang, J. Lai, Y. Li, Q. Lin, and J. S. Liu. Towards a statistical understanding of neural networks: Beyond the neural tangent kernel theories. *arXiv preprint arXiv:2412.18756*, 2024.