

# Global Convergence in Neural ODEs: Impact of Activation Functions

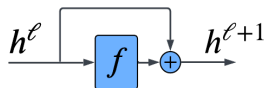
**Tianxiang Gao**<sup>1</sup>, Siyuan Sun<sup>2</sup>, Hailiang Liu<sup>2</sup>, Hongyang Gao<sup>2</sup>

<sup>1</sup>DePaul University

<sup>2</sup>Iowa State University

ICLR 2025, April 24, 2025

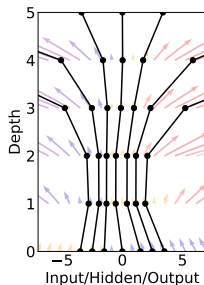
## Residual Networks (ResNet):



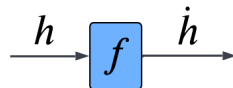
- Residual learning paradigm:

$$h_{\ell+1} = h_\ell + f(h_\ell, \theta_\ell)$$

- A discrete sequence of transformations:



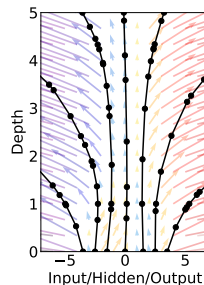
## Neural Ordinary Differential Equations (ODEs):



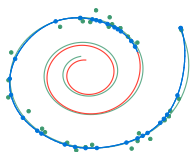
- Infinitely many layers and infinitesimal scaling:

$$\dot{h} = f(h_t, \theta),$$

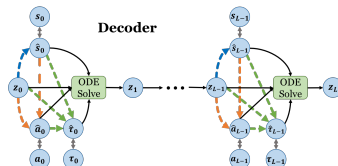
- A vector field continuously refining the state



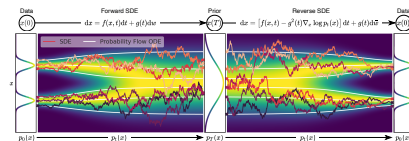
Thanks to their **continuous nature** and **parameter-sharing** efficiency, Neural ODEs have achieved success in a range of applications.



Time Series

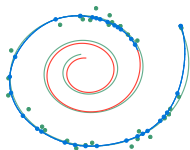


Reinforcement Learning

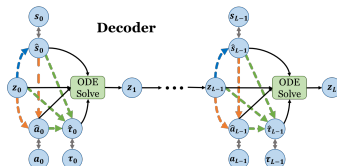


Score-based Generative Models

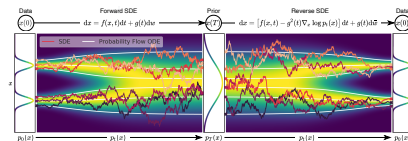
Thanks to their **continuous nature** and **parameter-sharing** efficiency, Neural ODEs have achieved success in a range of applications.



Time Series



Reinforcement Learning



Score-based Generative Models

## Key Question

As a **continuous** model with **shared** parameters, under what conditions does a Neural ODE **converge** under a gradient-based method?

- We study a Neural ODE  $f_{\theta}(x)$  defined by the **forward ODE** with hidden state  $\mathbf{h}_t$ :

$$\dot{\mathbf{h}}_t = \mathbf{W}\phi(\mathbf{h}_t), \quad \forall t \in [0, T] \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n)$ .

- We study a Neural ODE  $f_{\theta}(\mathbf{x})$  defined by the **forward ODE** with hidden state  $\mathbf{h}_t$ :

$$\dot{\mathbf{h}}_t = \mathbf{W}\phi(\mathbf{h}_t), \quad \forall t \in [0, T] \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n)$ .

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , we minimize the squared loss:

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} [f_{\theta}(\mathbf{x}_i) - y_i]^2 = \frac{1}{2} \|\mathbf{u} - \mathbf{y}\|^2, \quad (2)$$

where  $\mathbf{u}$  is the prediction vector with  $\mathbf{u}_i = f_{\theta}(\mathbf{x}_i)$ .

- We study a Neural ODE  $f_{\theta}(\mathbf{x})$  defined by the **forward ODE** with hidden state  $\mathbf{h}_t$ :

$$\dot{\mathbf{h}}_t = \mathbf{W}\phi(\mathbf{h}_t), \quad \forall t \in [0, T] \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n)$ .

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , we minimize the squared loss:

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} [f_{\theta}(\mathbf{x}_i) - y_i]^2 = \frac{1}{2} \|\mathbf{u} - \mathbf{y}\|^2, \quad (2)$$

where  $\mathbf{u}$  is the prediction vector with  $\mathbf{u}_i = f_{\theta}(\mathbf{x}_i)$ .

- Gradients are computed via the **backward ODE** of the adjoint state  $\lambda_t$ :

$$\dot{\lambda}_t = -\text{diag}(\phi'(\mathbf{h}_t))\mathbf{W}^{\top}, \quad \forall t \in [0, T] \quad (3)$$

- We study a Neural ODE  $f_{\theta}(x)$  defined by the **forward ODE** with hidden state  $\mathbf{h}_t$ :

$$\dot{\mathbf{h}}_t = \mathbf{W}\phi(\mathbf{h}_t), \quad \forall t \in [0, T] \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n)$ .

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , we minimize the squared loss:

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} [f_{\theta}(\mathbf{x}_i) - y_i]^2 = \frac{1}{2} \|\mathbf{u} - \mathbf{y}\|^2, \quad (2)$$

where  $\mathbf{u}$  is the prediction vector with  $\mathbf{u}_i = f_{\theta}(\mathbf{x}_i)$ .

- Gradients are computed via the **backward ODE** of the adjoint state  $\lambda_t$ :

$$\dot{\lambda}_t = -\text{diag}(\phi'(\mathbf{h}_t))\mathbf{W}^{\top}, \quad \forall t \in [0, T] \quad (3)$$

- Gradient with respect to weights:

$$\nabla_{\mathbf{W}} f_{\theta}(\mathbf{x}) = \int_0^T \lambda_t \phi(\mathbf{h}_t)^{\top} dt \quad (4)$$



- We study a Neural ODE  $f_{\theta}(x)$  defined by the **forward ODE** with hidden state  $\mathbf{h}_t$ :

$$\dot{\mathbf{h}}_t = \mathbf{W}\phi(\mathbf{h}_t), \quad \forall t \in [0, T] \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n)$ .

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , we minimize the squared loss:

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} [f_{\theta}(\mathbf{x}_i) - y_i]^2 = \frac{1}{2} \|\mathbf{u} - \mathbf{y}\|^2, \quad (2)$$

where  $\mathbf{u}$  is the prediction vector with  $\mathbf{u}_i = f_{\theta}(\mathbf{x}_i)$ .

- Gradients are computed via the **backward ODE** of the adjoint state  $\lambda_t$ :

$$\dot{\lambda}_t = -\text{diag}(\phi'(\mathbf{h}_t))\mathbf{W}^{\top}, \quad \forall t \in [0, T] \quad (3)$$

- Gradient with respect to weights:

$$\nabla_{\mathbf{W}} f_{\theta}(\mathbf{x}) = \int_0^T \lambda_t \phi(\mathbf{h}_t)^{\top} dt \quad (4)$$

- Parameters are updated by gradient descent:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} L(\theta^{(k)}) \quad (5)$$

## Well-Posedness and Gradient Computation

**Optimize-then-discretize:** Use numerical solvers to compute gradients by solving the continuous forward and backward ODEs.

**Optimize-then-discretize:** Use numerical solvers to compute gradients by solving the continuous forward and backward ODEs.

**Discretize-then-optimize:** Discretize the ODE using Euler's method, treating the system as a finite-depth ResNet  $f_{\theta}^L$  with shared parameters:

$$\mathbf{h}^{\ell} = \mathbf{h}^{\ell-1} + \Delta t \cdot \mathbf{W} \phi(\mathbf{h}^{\ell-1}), \quad \text{where } \Delta t = T/L \quad (6)$$

## Well-Posedness and Gradient Computation

**Optimize-then-discretize:** Use numerical solvers to compute gradients by solving the continuous forward and backward ODEs.

**Discretize-then-optimize:** Discretize the ODE using Euler's method, treating the system as a finite-depth ResNet  $f_{\theta}^L$  with shared parameters:

$$\mathbf{h}^{\ell} = \mathbf{h}^{\ell-1} + \Delta t \cdot \mathbf{W} \phi(\mathbf{h}^{\ell-1}), \quad \text{where } \Delta t = T/L \quad (6)$$

### Proposition 1

*Given  $T < \infty$ , if  $\phi$  is Lipschitz continuous, then the forward and backward ODEs have unique solutions  $\mathbf{h}_t$  and  $\boldsymbol{\lambda}_t$  for all  $t \in [0, T]$  and  $\mathbf{x} \in \mathbb{R}^d$ , almost surely over random initialization. Moreover,  $\boldsymbol{\lambda}_t(\mathbf{x}) = \partial f(\mathbf{x}; \boldsymbol{\theta}) / \partial \mathbf{h}_t$  is the solution to the backward ODE.*

## Well-Posedness and Gradient Computation

**Optimize-then-discretize:** Use numerical solvers to compute gradients by solving the continuous forward and backward ODEs.

**Discretize-then-optimize:** Discretize the ODE using Euler's method, treating the system as a finite-depth ResNet  $f_{\theta}^L$  with shared parameters:

$$\mathbf{h}^{\ell} = \mathbf{h}^{\ell-1} + \Delta t \cdot \mathbf{W} \phi(\mathbf{h}^{\ell-1}), \quad \text{where } \Delta t = T/L \quad (6)$$

### Proposition 1

Given  $T < \infty$ , if  $\phi$  is *Lipschitz continuous*, then the forward and backward ODEs have unique solutions  $\mathbf{h}_t$  and  $\lambda_t$  for all  $t \in [0, T]$  and  $\mathbf{x} \in \mathbb{R}^d$ , almost surely over random initialization. Moreover,  $\lambda_t(\mathbf{x}) = \partial f(\mathbf{x}; \theta) / \partial \mathbf{h}_t$  is the solution to the backward ODE.

### Proposition 2

If, in addition,  $\phi'$  is *Lipschitz continuous*, then the following holds a.s. over random initialization:

$$\left\| \nabla_{\theta} f^L(\mathbf{x}) - \nabla_{\theta} f(\mathbf{x}) \right\| \leq C L^{-1} \quad (7)$$

where  $C > 0$  is a constant depending on the Lipschitz constants, time horizon  $T$ , and  $\|\mathbf{x}\|$ .

- Previous studies have shown the training dynamics of predictions  $\mathbf{u}^k$  can be approximated by:

$$\mathbf{u}^{k+1} - \mathbf{y} \approx (\mathbf{I} - \eta \mathbf{H}^k)(\mathbf{u}^k - \mathbf{y}), \quad (8)$$

where  $\mathbf{H}^k \in \mathbb{R}^{N \times N}$  is the NTK Gram matrix:

$$K_{\theta}(\mathbf{x}, \bar{\mathbf{x}}) := \langle \nabla_{\theta} f(\mathbf{x}; \theta), \nabla_{\theta} f(\bar{\mathbf{x}}; \theta) \rangle \quad (9)$$

- Previous studies have shown the training dynamics of predictions  $\mathbf{u}^k$  can be approximated by:

$$\mathbf{u}^{k+1} - \mathbf{y} \approx (\mathbf{I} - \eta \mathbf{H}^k)(\mathbf{u}^k - \mathbf{y}), \quad (8)$$

where  $\mathbf{H}^k \in \mathbb{R}^{N \times N}$  is the NTK Gram matrix:

$$K_{\theta}(\mathbf{x}, \bar{\mathbf{x}}) := \langle \nabla_{\theta} f(\mathbf{x}; \theta), \nabla_{\theta} f(\bar{\mathbf{x}}; \theta) \rangle \quad (9)$$

- Neural ODE  $f_{\theta}$  can be viewed as the limit of a finite-depth ResNet  $f_{\theta}^L$  as  $L \rightarrow \infty$ .
- However, the NTK may depend on the order of limits:

$$\lim_{L \rightarrow \infty} \lim_{n \rightarrow \infty} f_{\theta}^L \neq \lim_{n \rightarrow \infty} \lim_{L \rightarrow \infty} f_{\theta}^L \quad (10)$$

- Previous studies have shown the training dynamics of predictions  $\mathbf{u}^k$  can be approximated by:

$$\mathbf{u}^{k+1} - \mathbf{y} \approx (\mathbf{I} - \eta \mathbf{H}^k)(\mathbf{u}^k - \mathbf{y}), \quad (8)$$

where  $\mathbf{H}^k \in \mathbb{R}^{N \times N}$  is the NTK Gram matrix:

$$K_{\theta}(\mathbf{x}, \bar{\mathbf{x}}) := \langle \nabla_{\theta} f(\mathbf{x}; \theta), \nabla_{\theta} f(\bar{\mathbf{x}}; \theta) \rangle \quad (9)$$

- Neural ODE  $f_{\theta}$  can be viewed as the limit of a finite-depth ResNet  $f_{\theta}^L$  as  $L \rightarrow \infty$ .
- However, the NTK may depend on the order of limits:

$$\lim_{L \rightarrow \infty} \lim_{n \rightarrow \infty} f_{\theta}^L \neq \lim_{n \rightarrow \infty} \lim_{L \rightarrow \infty} f_{\theta}^L \quad (10)$$

- Illustrative example:

$$a_{n,\ell} := \frac{n}{\ell + n}, \quad \lim_{\ell \rightarrow \infty} \lim_{n \rightarrow \infty} a_{n,\ell} = 1, \quad \lim_{n \rightarrow \infty} \lim_{\ell \rightarrow \infty} a_{n,\ell} = 0 \quad (11)$$



- Previous studies have shown the training dynamics of predictions  $\mathbf{u}^k$  can be approximated by:

$$\mathbf{u}^{k+1} - \mathbf{y} \approx (\mathbf{I} - \eta \mathbf{H}^k)(\mathbf{u}^k - \mathbf{y}), \quad (8)$$

where  $\mathbf{H}^k \in \mathbb{R}^{N \times N}$  is the NTK Gram matrix:

$$K_{\theta}(\mathbf{x}, \bar{\mathbf{x}}) := \langle \nabla_{\theta} f(\mathbf{x}; \theta), \nabla_{\theta} f(\bar{\mathbf{x}}; \theta) \rangle \quad (9)$$

- Neural ODE  $f_{\theta}$  can be viewed as the limit of a finite-depth ResNet  $f_{\theta}^L$  as  $L \rightarrow \infty$ .
- However, the NTK may depend on the order of limits:

$$\lim_{L \rightarrow \infty} \lim_{n \rightarrow \infty} f_{\theta}^L \neq \lim_{n \rightarrow \infty} \lim_{L \rightarrow \infty} f_{\theta}^L \quad (10)$$

- Illustrative example:

$$a_{n,\ell} := \frac{n}{\ell + n}, \quad \lim_{\ell \rightarrow \infty} \lim_{n \rightarrow \infty} a_{n,\ell} = 1, \quad \lim_{n \rightarrow \infty} \lim_{\ell \rightarrow \infty} a_{n,\ell} = 0 \quad (11)$$

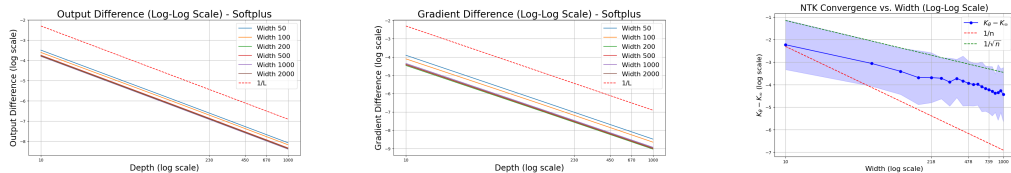
## Theorem 1 (NTK Convergence)

Suppose  $\phi$  and  $\phi'$  are Lipschitz. Then as width  $n \rightarrow \infty$ , the NTK  $K_{\theta}$  of the Neural ODE  $f_{\theta}$  converges almost surely to a deterministic kernel:

$$K_{\theta} \rightarrow K_{\infty}, \quad \text{as } n \rightarrow \infty, \quad (12)$$

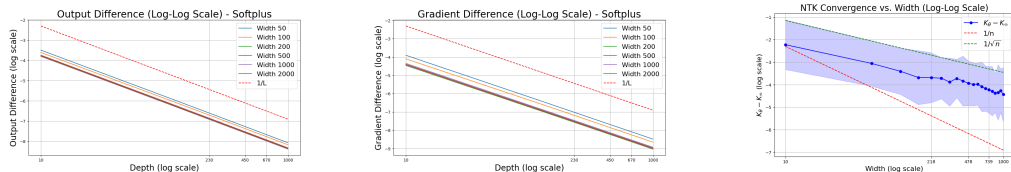
where  $K_{\infty}$  equals the limit of NTKs  $K_{\infty}^L$  for ResNets as depth  $L \rightarrow \infty$ .

# Empirical Validation of Well-Posedness and NTK Convergence

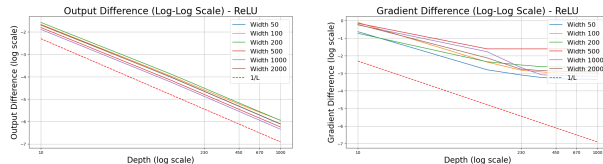


**Figure:** Comparison of Neural ODEs and ResNets: With Softplus (smooth), output, gradient, and NTK differences decay as  $1/L$ .

# Empirical Validation of Well-Posedness and NTK Convergence



**Figure:** Comparison of Neural ODEs and ResNets: With Softplus (smooth), output, gradient, and NTK differences decay as  $1/L$ .



**Figure:** Comparison of Neural ODEs and ResNets with **ReLU (non-smooth)** activations: Output difference shows  $1/L$  convergence, while gradient difference remains constant as depth  $L$  increases.

- Previous works have shown a strictly positive definite (SPD) limiting NTK  $K_\infty$  is sufficient to guarantee global convergence under gradient descent.

- Previous works have shown a strictly positive definite (SPD) limiting NTK  $K_\infty$  is sufficient to guarantee global convergence under gradient descent.
- However, the SPD property of the finite-depth NTK  $K_\infty^L$  does not ensure that  $K_\infty$  remains SPD as  $L \rightarrow \infty$ . The spectrum may degenerate with increasing depth.

- Previous works have shown a strictly positive definite (SPD) limiting NTK  $K_\infty$  is sufficient to guarantee global convergence under gradient descent.
- However, the SPD property of the finite-depth NTK  $K_\infty^L$  does not ensure that  $K_\infty$  remains SPD as  $L \rightarrow \infty$ . The spectrum may degenerate with increasing depth.
- The NTK of the Neural ODE admits an **integral form**:

$$K_\infty = \Sigma^{T,T} + \int_0^T \int_0^T \Sigma^{t,s} K^{t,s} dt ds + \Sigma^{0,0} K^{0,0}$$

where  $\Sigma^{t,s}$  and  $K^{t,s}$  encode covariance and kernel terms at time  $t, s$ .

- Previous works have shown a strictly positive definite (SPD) limiting NTK  $K_\infty$  is sufficient to guarantee global convergence under gradient descent.
- However, the SPD property of the finite-depth NTK  $K_\infty^L$  does not ensure that  $K_\infty$  remains SPD as  $L \rightarrow \infty$ . The spectrum may degenerate with increasing depth.
- The NTK of the Neural ODE admits an **integral form**:

$$K_\infty = \Sigma^{T,T} + \int_0^T \int_0^T \Sigma^{t,s} K^{t,s} dt ds + \Sigma^{0,0} K^{0,0}$$

where  $\Sigma^{t,s}$  and  $K^{t,s}$  encode covariance and kernel terms at time  $t, s$ .

### Proposition 3

*Suppose  $\phi$  and  $\phi'$  are Lipschitz continuous. If  $\phi$  is nonlinear and non-polynomial, then the limiting NTK  $K_\infty$  is strictly positive definite.*

## Assumption 1

Let  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  be a training set. Assume

- ① *Training set:*  $\mathbf{x}_i \in \mathbb{S}^{d-1}$  and  $\mathbf{x}_i \neq \mathbf{x}_j$  for all  $i \neq j$ ;  $|y_i| = \mathcal{O}(1)$ ,
- ② *Smoothness:*  $\phi$  and  $\phi'$  are Lipschitz continuous, respectively,
- ③ *Nonlinearity:*  $\phi$  is nonlinear and non-polynomial.



## Assumption 1

Let  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  be a training set. Assume

- ① *Training set:*  $\mathbf{x}_i \in \mathbb{S}^{d-1}$  and  $\mathbf{x}_i \neq \mathbf{x}_j$  for all  $i \neq j$ ;  $|y_i| = \mathcal{O}(1)$ ,
- ② *Smoothness:*  $\phi$  and  $\phi'$  are Lipschitz continuous, respectively,
- ③ *Nonlinearity:*  $\phi$  is nonlinear and non-polynomial.

## Theorem 2

Suppose Assumption 1 holds and the learning rate  $\eta$  is chosen s.t.  $0 < \eta \leq 1/\|\mathbf{X}\|^2$ . Then for any  $\delta > 0$ , there exists a natural number  $n_\delta$  s.t. for all widths  $n \geq n_\delta$  the following results hold with probability at least  $1 - \delta$  over random initialization:

- ① The parameters  $\boldsymbol{\theta}^k$  stay in a neighborhood of  $\boldsymbol{\theta}^0$ , i.e.,

$$\|\boldsymbol{\theta}^k - \boldsymbol{\theta}^0\| \leq C\|\mathbf{X}\|\sqrt{L(\boldsymbol{\theta}_0)}/\lambda_0. \quad (13)$$

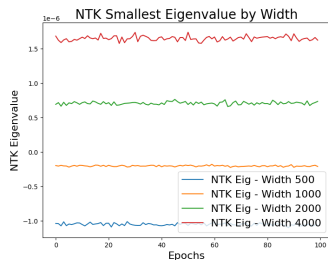
- ② The loss function  $L(\boldsymbol{\theta}_k)$  consistently decreases to zero at an exponential rate, i.e.,

$$L(\boldsymbol{\theta}_k) \leq (1 - \eta\lambda_0)^k L(\boldsymbol{\theta}_0), \quad (14)$$

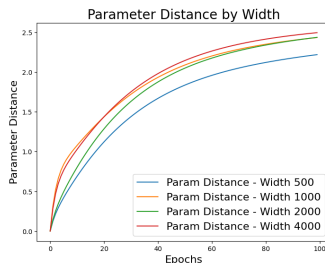
where  $\lambda_0 := \lambda_{\min}(K_\infty) > 0$ , and the constant  $C > 0$  only depends on Lipschitz coefficients and  $T$ .

# Empirical Validation of Global Convergence

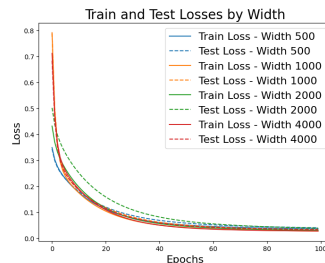
- We examine NTK spectra, parameter stability, and loss convergence across different model widths.



(a) NTK eigenvalues



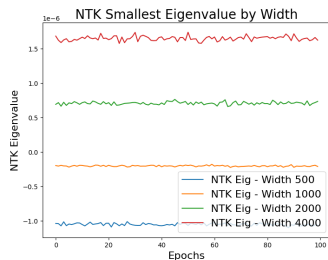
(b) Distance from initialization



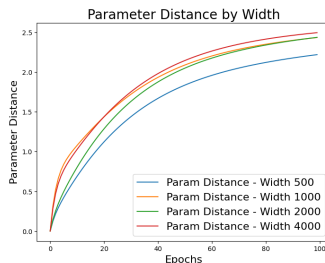
(c) Train vs. test loss

# Empirical Validation of Global Convergence

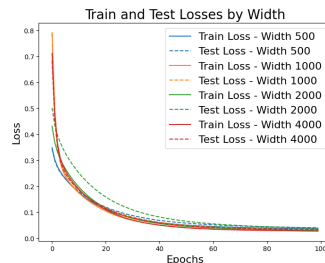
- We examine NTK spectra, parameter stability, and loss convergence across different model widths.



(a) NTK eigenvalues



(b) Distance from initialization



(c) Train vs. test loss

- Additional experiments: time horizon effects, quadratic activations, solver comparisons (adaptive vs. fixed step). See the **poster** or **paper** for details.

### Open Directions

- **Convergence Rates:** Derive explicit rates as width/depth grow in Neural ODEs.
- **Solver Design:** Design efficient ODE solvers guided by gradient alignment conditions.
- **Beyond GD:** Explore SGD, momentum, and adaptive optimization methods.
- **Feature Learning:** Analyze training and generalization beyond the lazy training regime.
- **Broader Models:** Extend theory to Transformers, ResNets, and state-space models.

### Contact Information:

- Email: [tgao9@depaul.edu](mailto:tgao9@depaul.edu)
- Web: <https://gaotx-cs.github.io/>