

# THE SURPRISING SIMPLICITY OF OVERPARAMETERIZED DEEP NEURAL NETWORKS

---

JEFFREY PENNINGTON

GOOGLE BRAIN

STATS 385

10-16-19

# COLLABORATION

Jascha Sohl-Dickstein

Sam Schoenholz

Surya Ganguli

Greg Yang

Lechao Xiao

Yasaman Bahri

Jaehoon Lee

Roman Novak

Minmin Chen

Dar Gilboa

Bo Chang

# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion

# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion



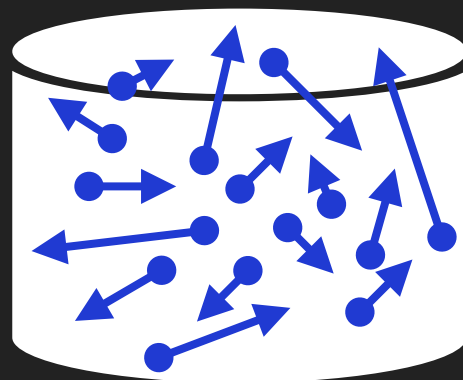
MOTIVATION: WHY STUDY OVERPARAMETERIZED MODELS?

---

# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

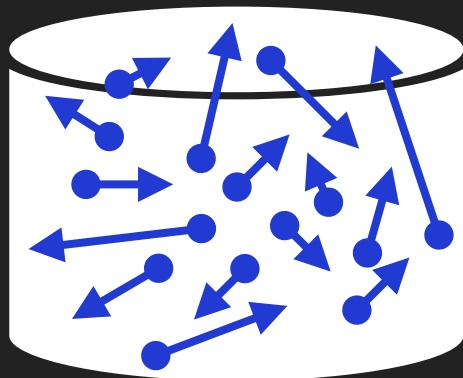
# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

Microscopic



# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

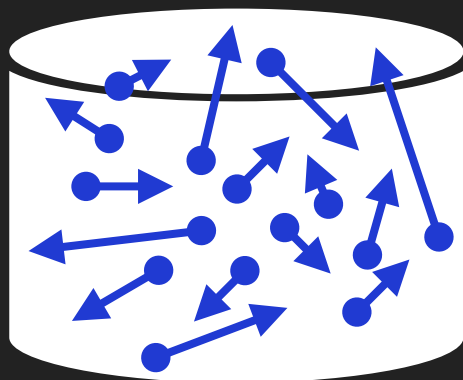
Microscopic



$$\{p_i, x_i\}_{i=1\dots N}$$

# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

Microscopic

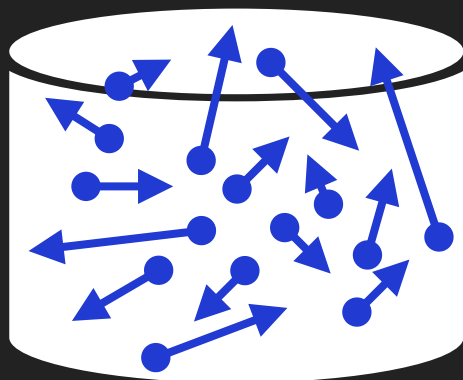


$$\{p_i, x_i\}_{i=1\dots N}$$

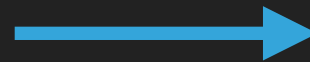
$$H = \frac{1}{2m} \sum_{i=1}^N p_i^2 + \sum_{i=1}^N V(x_i) + \sum_{i < j} U(x_i - x_j)$$

# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

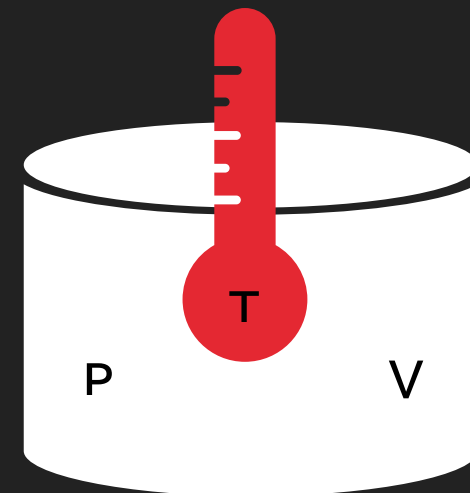
Microscopic



$N \gg 1$



Macroscopic

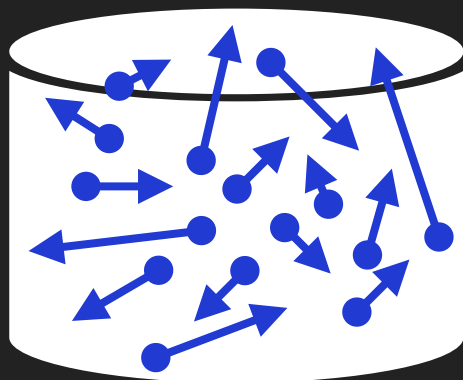


$$\{p_i, x_i\}_{i=1 \dots N}$$

$$H = \frac{1}{2m} \sum_{i=1}^N p_i^2 + \sum_{i=1}^N V(x_i) + \sum_{i < j} U(x_i - x_j)$$

# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

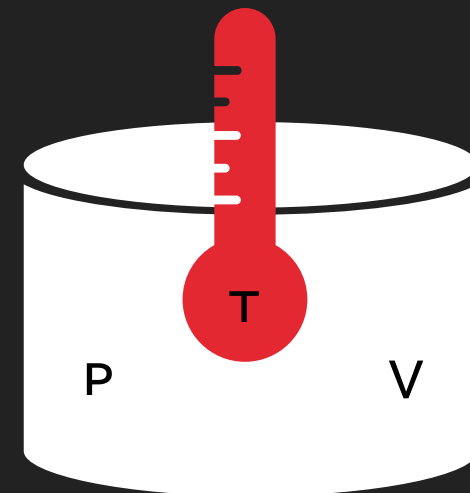
Microscopic



$$\{p_i, x_i\}_{i=1\dots N}$$



Macroscopic

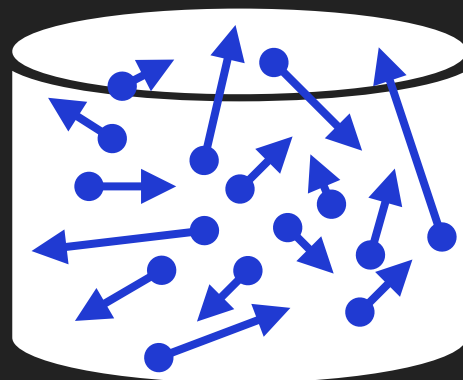


$$\{P, V, T\}$$

$$H = \frac{1}{2m} \sum_{i=1}^N p_i^2 + \sum_{i=1}^N V(x_i) + \sum_{i<j} U(x_i - x_j)$$

# SIMPLICITY IN LARGE NUMBERS: STATISTICAL MECHANICS

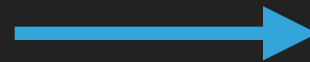
Microscopic



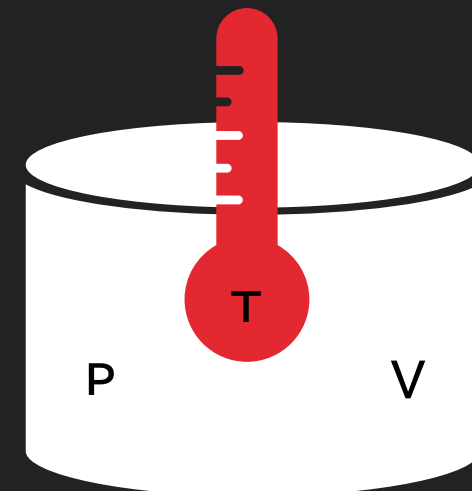
$$\{p_i, x_i\}_{i=1\dots N}$$

$$H = \frac{1}{2m} \sum_{i=1}^N p_i^2 + \sum_{i=1}^N V(x_i) + \sum_{i<j} U(x_i - x_j)$$

$N \gg 1$



Macroscopic



$$\{P, V, T\}$$

$$PV = nRT$$

MOTIVATION: WHY STUDY OVERPARAMETERIZED MODELS?

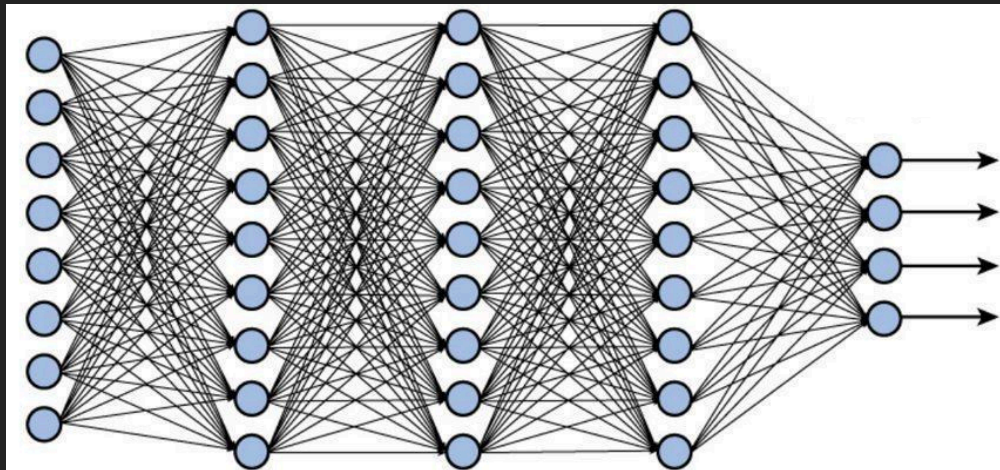
---

# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS



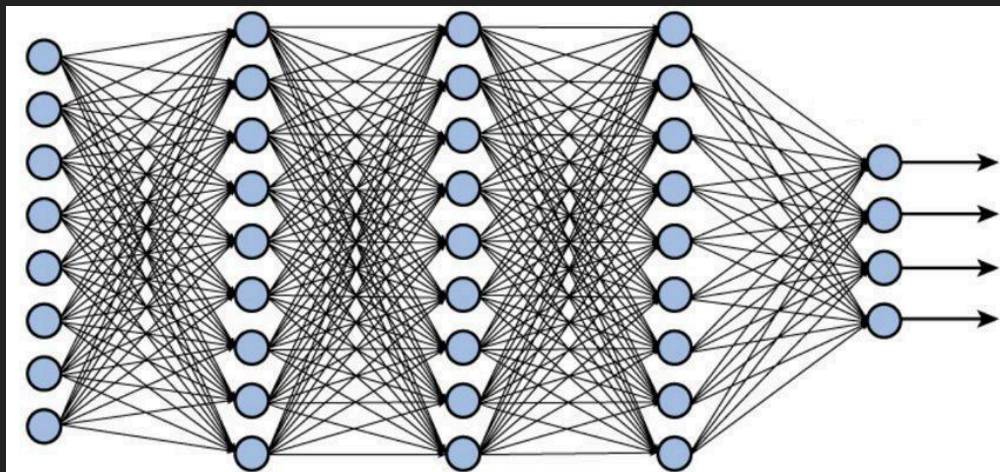
# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

Microscopic



# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

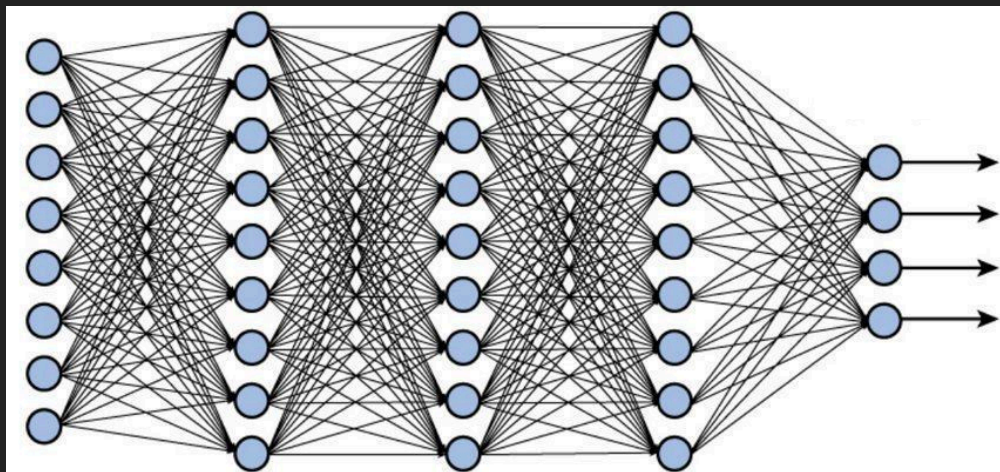
Microscopic



$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

Microscopic

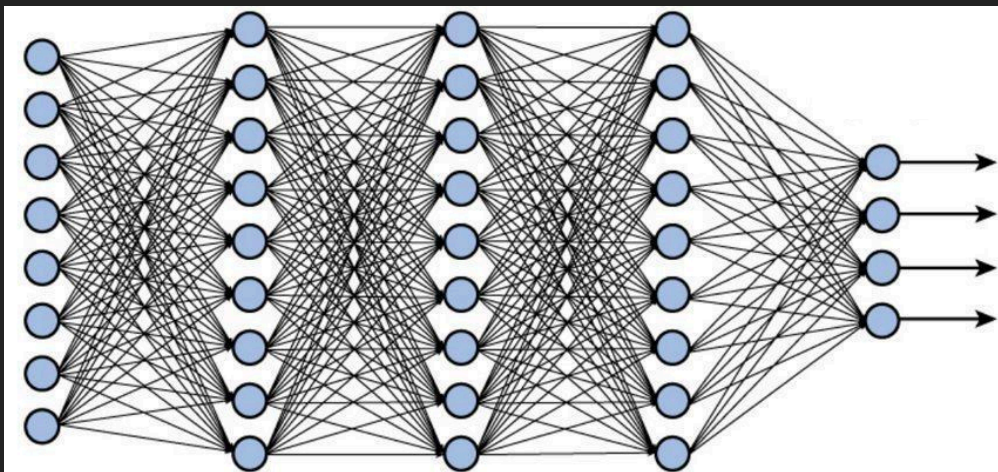


$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

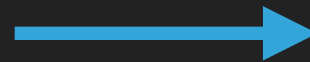
$$\partial_t \Theta = -\nabla_{\Theta} \mathcal{L}$$

# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

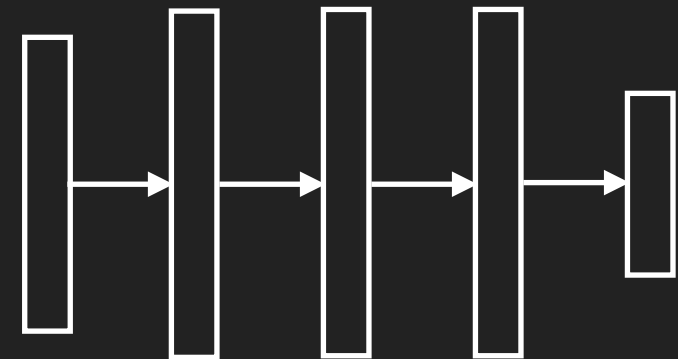
Microscopic



$N \gg 1$



Macroscopic



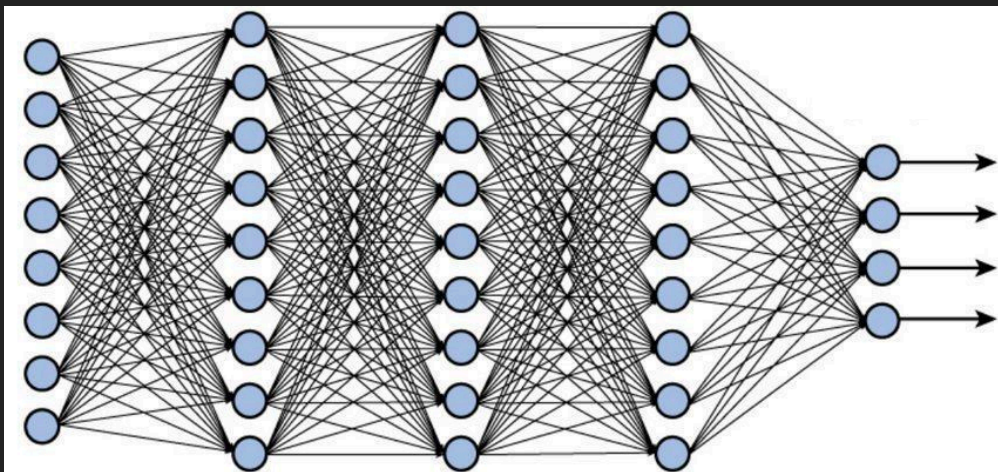
$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

$$\partial_t \Theta = -\nabla_{\Theta} \mathcal{L}$$

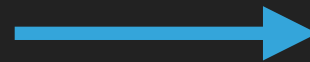


# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

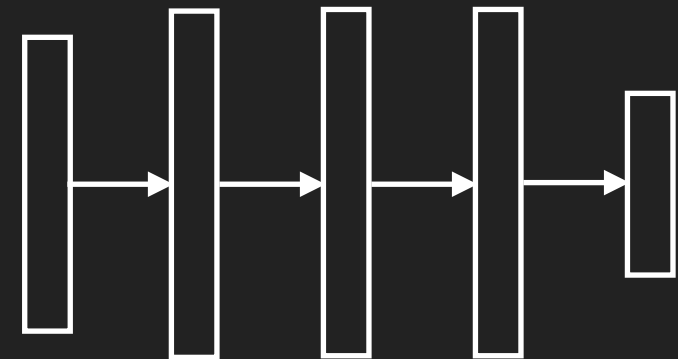
Microscopic



$N \gg 1$



Macroscopic



$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

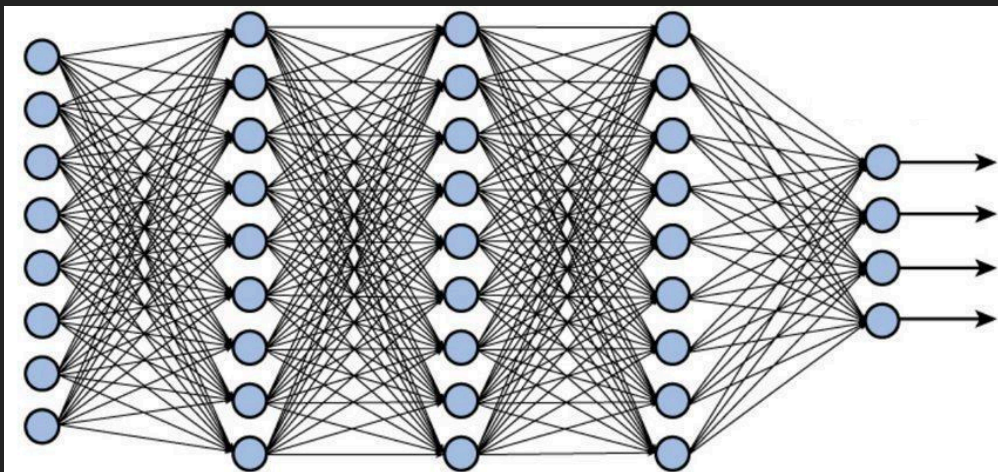
?

$$\partial_t \Theta = -\nabla_{\Theta} \mathcal{L}$$

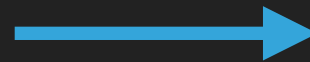
?

# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

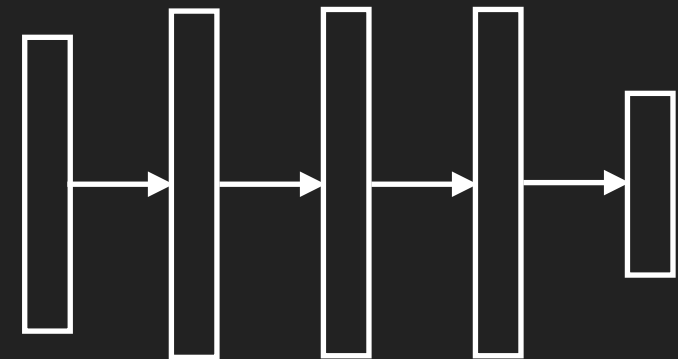
Microscopic



$N \gg 1$



Macroscopic



$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

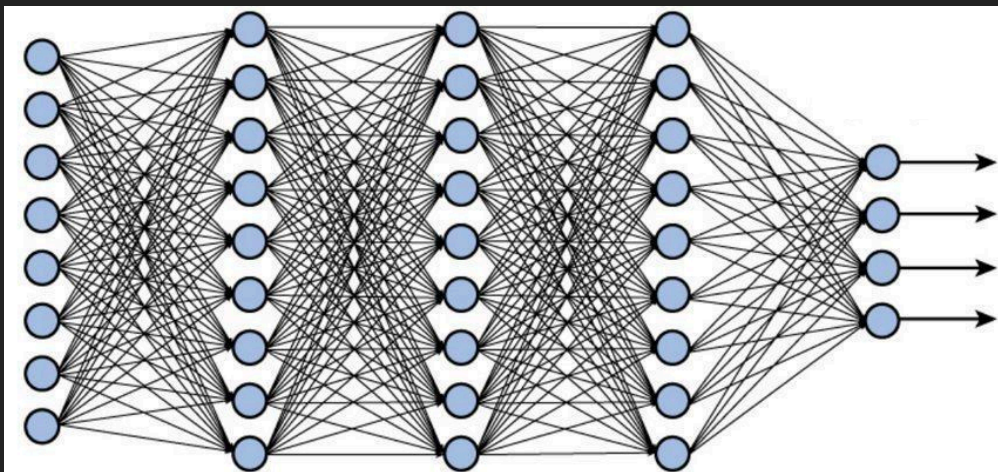
$$\Sigma_{ab}^l = \frac{1}{N} \sum_{i=1}^N x_{ia}^l x_{ib}^l \quad (?)$$

$$\partial_t \Theta = -\nabla_{\Theta} \mathcal{L}$$

?

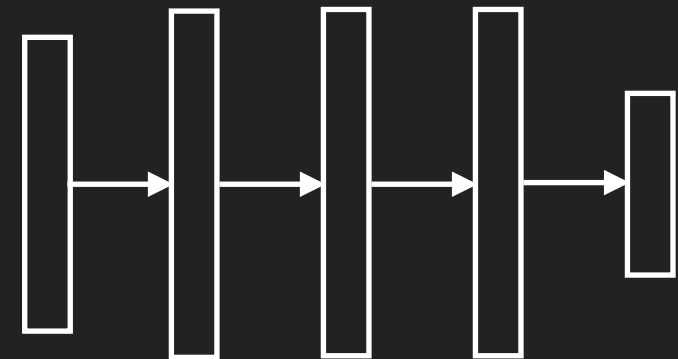
# SIMPLICITY IN LARGE NUMBERS: NEURAL NETWORKS

Microscopic



$N \gg 1$   
→

Macroscopic



$$\Theta = \{W_{ij}^l, b_i^l\}_{i,j=1\dots N}^{l=1\dots L}$$

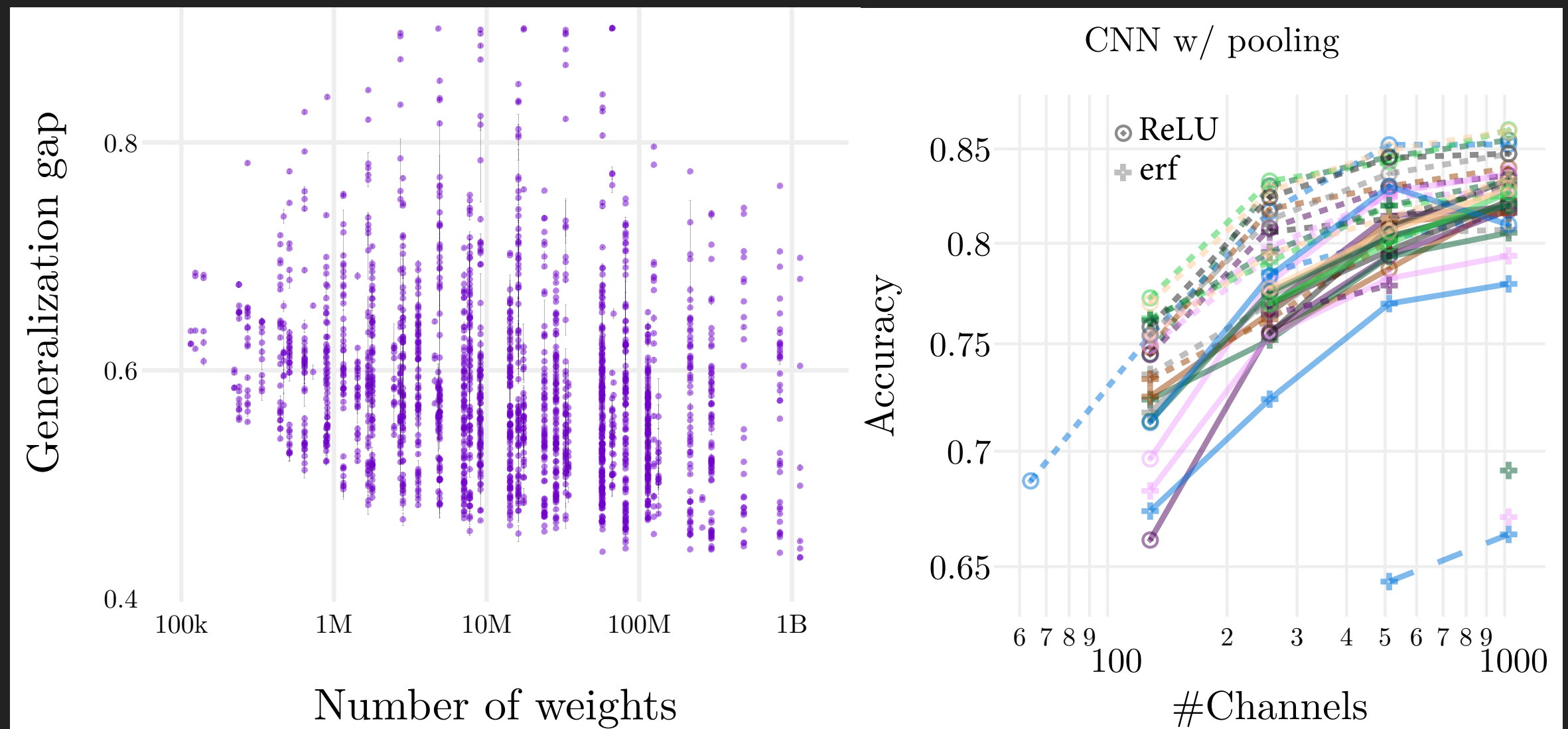
$$\partial_t \Theta = -\nabla_{\Theta} \mathcal{L}$$

$$\Sigma_{ab}^l = \frac{1}{N} \sum_{i=1}^N x_{ia}^l x_{ib}^l \quad (?)$$

$$\partial_t \Sigma_{ab}^l \approx 0 \quad (?)$$

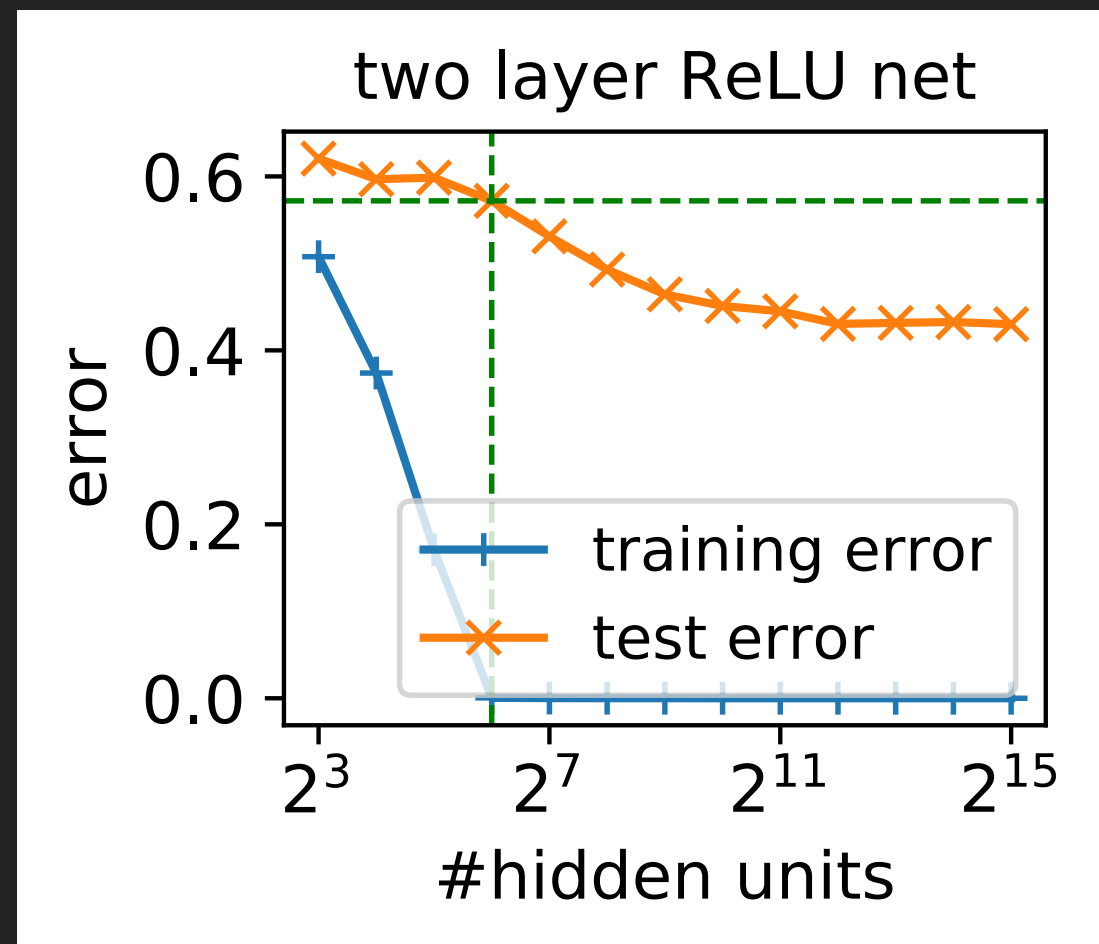
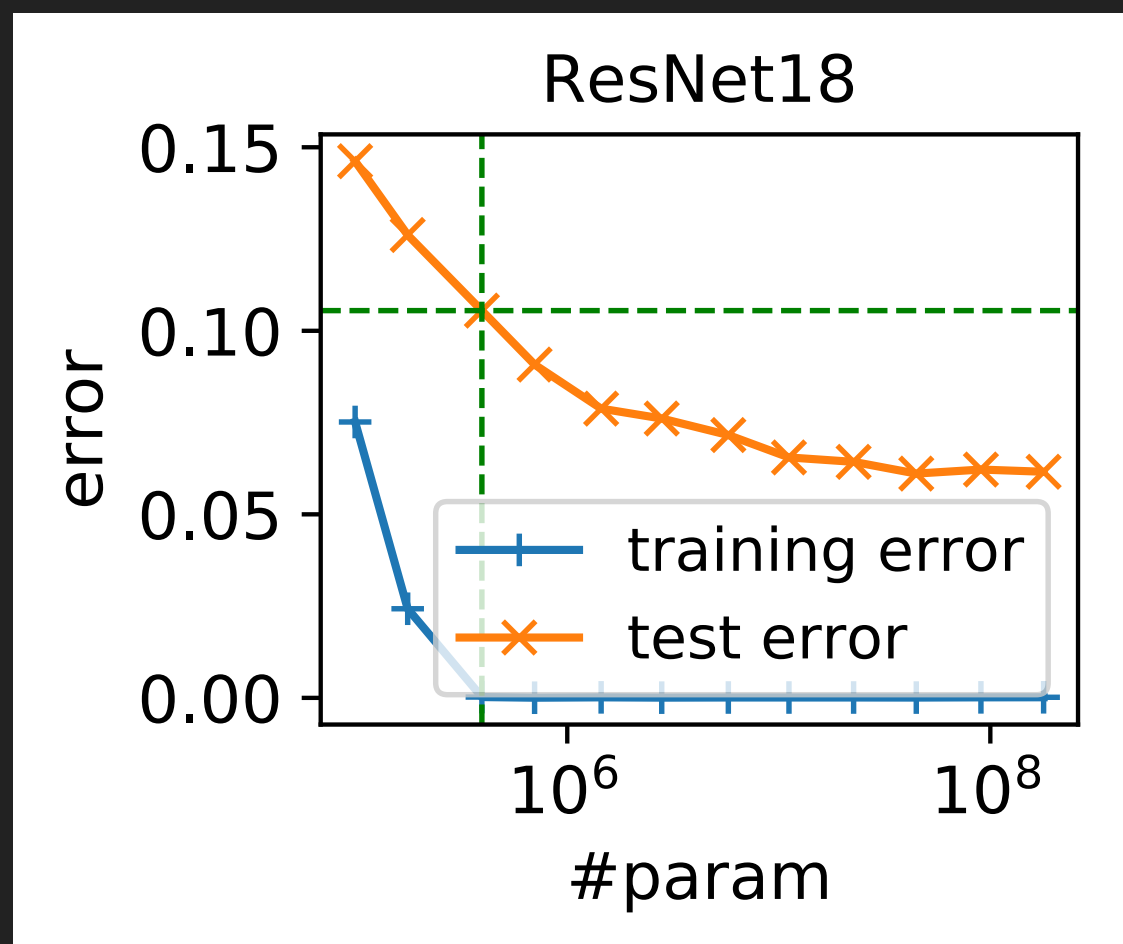
## MOTIVATION: WHY STUDY OVERPARAMETERIZED MODELS?

# OVERPARAMETERIZED MODELS PERFORM BETTER





# OVERPARAMETERIZED MODELS PERFORM BETTER



# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion

## THE SINGLE HIDDEN LAYER CASE

Fully connected, single hidden layer [Radford Neal '94]

Inputs:  $\mathbf{x}_a \in \mathbb{R}^{N_0}$  with input index  $a$

$$\Sigma_{ab}^0 = \frac{1}{N_0} \sum_i x_{ia} x_{ib}$$

Parameters:  $W_{ij}^l \in \mathbb{R}^{N_{l-1} \times N_l}$      $b_i^l \in \mathbb{R}^{N_l}$

Prior:  $W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2 / N_{l-1})$      $b_i^l \sim \mathcal{N}(0, \sigma_b^2)$

Network:

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1$$
$$y_{ia} = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

## THE SINGLE HIDDEN LAYER CASE

Fully connected, single hidden layer [Radford Neal '94]

Inputs:  $\mathbf{x}_a \in \mathbb{R}^{N_0}$  with input index  $a$

$$\Sigma_{ab}^0 = \frac{1}{N_0} \sum_i x_{ia} x_{ib}$$

Parameters:  $W_{ij}^l \in \mathbb{R}^{N_{l-1} \times N_l}$      $b_i^l \in \mathbb{R}^{N_l}$

Prior:  $W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2 / N_{l-1})$      $b_i^l \sim \mathcal{N}(0, \sigma_b^2)$

Network:

Weighted sum of Gaussians

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1$$

$$(z_{ia}^1, z_{jb}^1)^T \sim \mathcal{N}(0, \Sigma_{ab}^1 \delta_{ij})$$

$$y_{ia} = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

## THE SINGLE HIDDEN LAYER CASE

Fully connected, single hidden layer [Radford Neal '94]

Inputs:  $\mathbf{x}_a \in \mathbb{R}^{N_0}$  with input index  $a$

$$\Sigma_{ab}^0 = \frac{1}{N_0} \sum_i x_{ia} x_{ib}$$

Parameters:  $W_{ij}^l \in \mathbb{R}^{N_{l-1} \times N_l}$      $b_i^l \in \mathbb{R}^{N_l}$

Prior:  $W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2 / N_{l-1})$      $b_i^l \sim \mathcal{N}(0, \sigma_b^2)$

Network:

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1$$

Weighted sum of Gaussians

$$(z_{ia}^1, z_{jb}^1)^T \sim \mathcal{N}(0, \Sigma_{ab}^1 \delta_{ij})$$

$$y_{ia} = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

Sum of i.i.d. random variables

$$(y_{ia}, y_{jb})^T \xrightarrow{N_1 \rightarrow \infty} \mathcal{N}(0, \Sigma_{ab}^2 \delta_{ij})$$

## THE SINGLE HIDDEN LAYER CASE

Infinitely wide neural networks are Gaussian Processes

$$\begin{aligned} z_{ia}^1 &= \sum_j W_{ij}^1 x_{ja} + b_i^1 & (z_{ia}^1, z_{jb}^1)^T &\sim \mathcal{N}(0, \Sigma_{ab}^1 \delta_{ij}) \\ y_{ia} &= \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2 & (y_{ia}, y_{jb})^T &\xrightarrow{N_1 \rightarrow \infty} \mathcal{N}(0, \Sigma_{ab}^2 \delta_{ij}) \end{aligned}$$

Completely defined by a compositional kernel

## THE SINGLE HIDDEN LAYER CASE

Infinitely wide neural networks are Gaussian Processes

$$\begin{aligned} z_{ia}^1 &= \sum_j W_{ij}^1 x_{ja} + b_i^1 & (z_{ia}^1, z_{jb}^1)^T &\sim \mathcal{N}(0, \Sigma_{ab}^1 \delta_{ij}) \\ y_{ia} &= \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2 & (y_{ia}, y_{jb})^T &\xrightarrow{N_1 \rightarrow \infty} \mathcal{N}(0, \Sigma_{ab}^2 \delta_{ij}) \end{aligned}$$

Completely defined by a compositional kernel

$$\Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$

$$\Sigma^2 = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, \Sigma^1)} [\phi(\mathbf{z}) \phi(\mathbf{z})^T] + \sigma_b^2$$

Significant simplification

# THE SINGLE HIDDEN LAYER CASE

Infinitely wide neural networks are Gaussian Processes

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1$$

$$(z_{ia}^1, z_{jb}^1)^T \sim \mathcal{N}(0, \Sigma_{ab}^1 \delta_{ij})$$

$$y_{ia} = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

$$(y_{ia}, y_{jb})^T \xrightarrow{N_1 \rightarrow \infty} \mathcal{N}(0, \Sigma_{ab}^2 \delta_{ij})$$

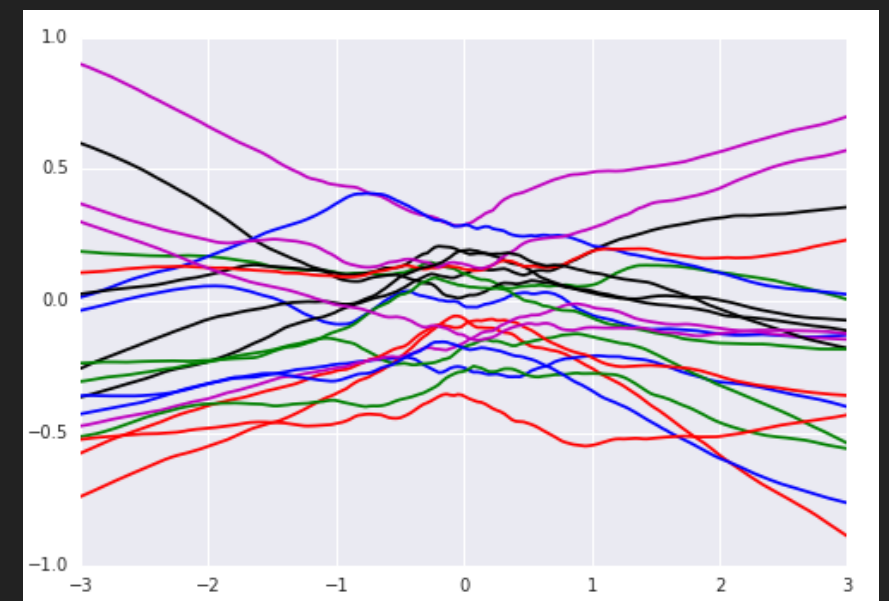
Completely defined by a compositional kernel

$$\Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$

$$\Sigma^2 = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^1)} [\phi(z) \phi(z)^T] + \sigma_b^2$$

Significant simplification

Draws from ReLU-GP





## DEEP NETWORKS

Extension to deep networks

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1 \longrightarrow \Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$

[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)


[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)

## DEEP NETWORKS

Extension to deep networks

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1 \quad \longrightarrow \quad \Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$



$$z_{ia}^2 = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)

## DEEP NETWORKS

Extension to deep networks

[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1 \longrightarrow \Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$

$$\downarrow$$
$$z_{ia}^2 = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2 \xrightarrow[N_1 \rightarrow \infty]{\text{---}} \Sigma^2 = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \Sigma^1)} [\phi(\mathbf{z}) \phi(\mathbf{z})^T] + \sigma_b^2$$

## DEEP NETWORKS

Extension to deep networks

[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)

$$z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1$$



$$\Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2$$



$$z_{ia}^2 = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2$$

$N_1 \rightarrow \infty$



$$\Sigma^2 = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \Sigma^1)} [\phi(\mathbf{z}) \phi(\mathbf{z})^T] + \sigma_b^2$$

## DEEP NETWORKS

Extension to deep networks

[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)

$$\begin{array}{ccc}
 z_{ia}^1 = \sum_j W_{ij}^1 x_{ja} + b_i^1 & \xrightarrow{\quad} & \Sigma^1 = \sigma_w^2 \Sigma^0 + \sigma_b^2 \\
 \downarrow & & \downarrow \mathcal{C} \\
 z_{ia}^2 = \sum_j W_{ij}^2 \phi(z_{ja}^1) + b_i^2 & \xrightarrow[N_1 \rightarrow \infty]{\quad} & \Sigma^2 = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, \Sigma^1)} [\phi(\mathbf{z}) \phi(\mathbf{z})^T] + \sigma_b^2 \\
 \downarrow & & \downarrow \mathcal{C} \\
 \vdots & & \vdots \\
 \downarrow & & \downarrow \mathcal{C} \\
 z_{ia}^l = \sum_j W_{ij}^l \phi(z_{ja}^{l-1}) + b_i^l & \xrightarrow[N_{l-1} \rightarrow \infty]{\quad} & \Sigma^l = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, \Sigma^{l-1})} [\phi(\mathbf{z}) \phi(\mathbf{z})^T] + \sigma_b^2
 \end{array}$$

## DEEP NETWORKS

Extension to deep networks

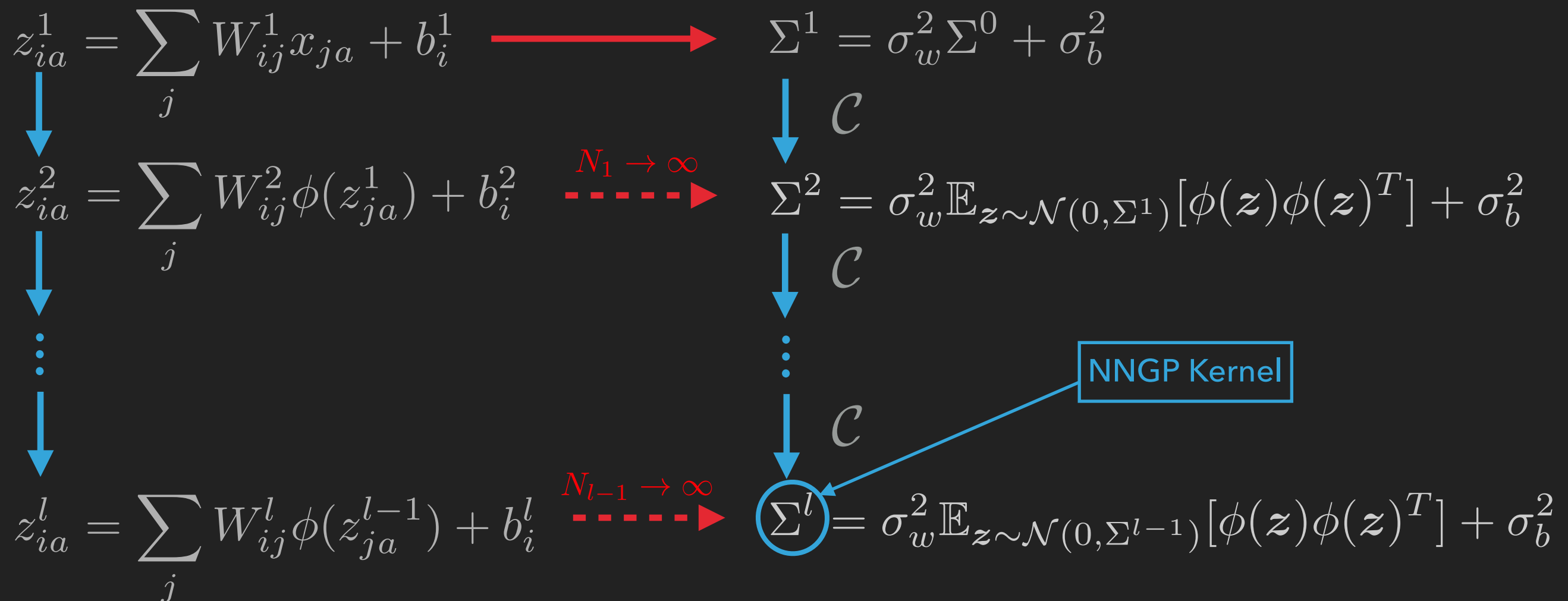
[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)



## DEEP NETWORKS

Extension to deep networks

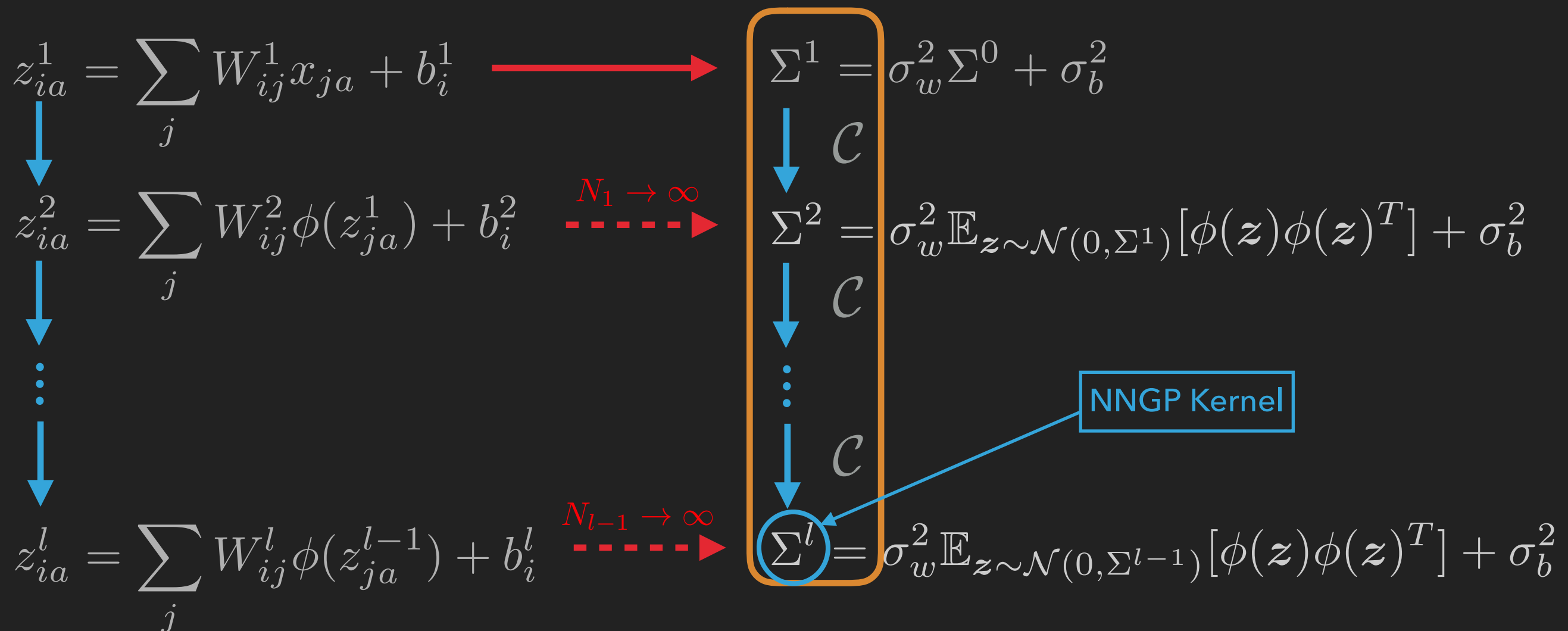
[AD, RF, YS \('16\)](#)

[BP, SL, MR, JSD, SG \('16\)](#)

[SSS, JG, SG, JSD \('17\)](#)

[JL, YB et al. \('18\)](#)

[RN, XLC et al. \('18\)](#)



Neural network induces **dynamical system** over kernels

Understanding prior equivalent to studying dynamics

# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion



# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots$$

# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

Rate of convergence determined by behavior near fixed point

# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

Rate of convergence determined by behavior near fixed point

$$\epsilon^l = \Sigma^* - \Sigma^l \quad \Rightarrow \quad \epsilon^{l+1} = \left. \frac{\partial \mathcal{C}(\Sigma)}{\partial \Sigma} \right|_{\Sigma^*} \epsilon^l$$

# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

Rate of convergence determined by behavior near fixed point

$$\epsilon^l = \Sigma^* - \Sigma^l \quad \Rightarrow \quad \epsilon^{l+1} = \left. \frac{\partial \mathcal{C}(\Sigma)}{\partial \Sigma} \right|_{\Sigma^*} \epsilon^l$$

$\lambda_{\max} > 1$   
 $\swarrow$   
 Unstable fixed point

$\searrow$   
 $\lambda_{\max} \leq 1$   
 Stable fixed point

# DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

Rate of convergence determined by behavior near fixed point

$$\epsilon^l = \Sigma^* - \Sigma^l \quad \Rightarrow \quad \epsilon^{l+1} \approx \lambda_{\max}^l$$

$$\lambda_{\max} > 1$$

Unstable fixed point

$$\lambda_{\max} \leq 1$$

Stable fixed point

## DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

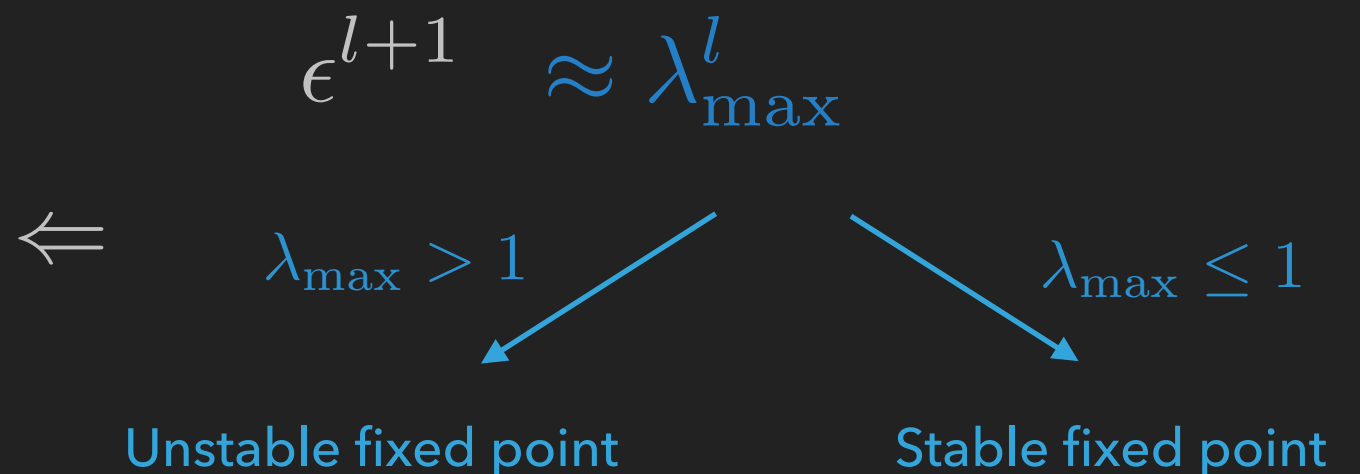
Rate of convergence determined by behavior near fixed point

Exponential convergence

$$\epsilon^l \approx e^{-l/\xi}$$

with rate

$$\xi = -1/\log \lambda_{\max}$$



## DYNAMICS OF SIGNAL PROPAGATION

$$\Sigma^1 \xrightarrow{\mathcal{C}} \Sigma^2 \xrightarrow{\mathcal{C}} \dots \xrightarrow{\mathcal{C}} \Sigma^l \xrightarrow{\mathcal{C}} \dots \Sigma^*$$

Dynamics converge to **universal** fixed point

- Independent of inputs  $\Rightarrow$  pathological

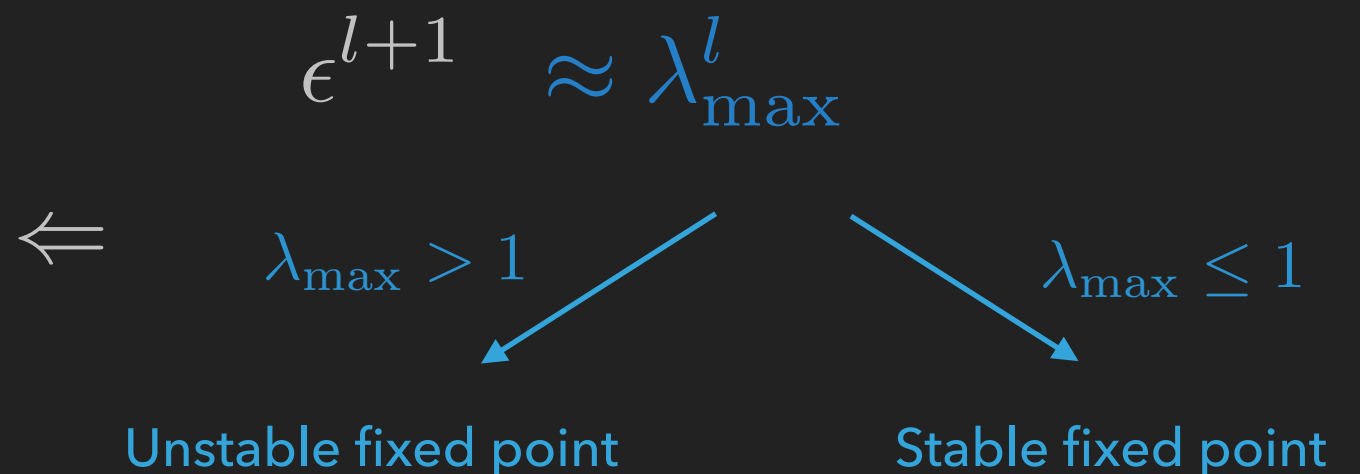
Rate of convergence determined by behavior near fixed point

Exponential convergence

$$\epsilon^l \approx e^{-l/\xi}$$

with rate

$$\xi = -1/\log \lambda_{\max}$$



How can we adjust the hyperparameters to delay convergence?



## FIXED POINT ANALYSIS

The fixed point satisfies


$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

## FIXED POINT ANALYSIS

The fixed point satisfies

$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

One solution is perfect correlation,

$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$



Variance      Correlation

## FIXED POINT ANALYSIS

The fixed point satisfies

$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

One solution is perfect correlation,

$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$


Variance      Correlation


Is this fixed point stable?

## FIXED POINT ANALYSIS

The fixed point satisfies

$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

One solution is perfect correlation,

$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$


Variance      Correlation

Is this fixed point stable?

- Depends on hyperparameters

# FIXED POINT ANALYSIS

The fixed point satisfies

$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

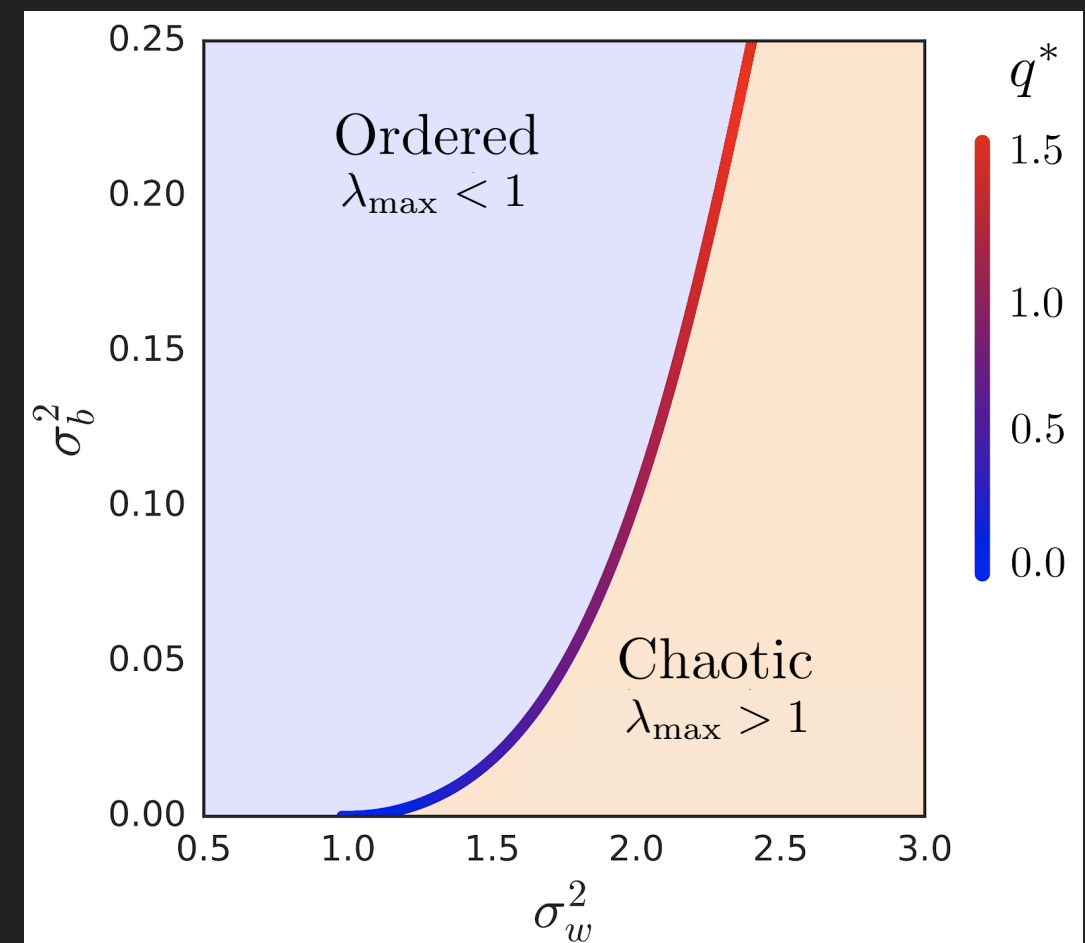
One solution is perfect correlation,

$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

↑ Variance
 ↙ Correlation

Is this fixed point stable?

- Depends on hyperparameters



# FIXED POINT ANALYSIS

The fixed point satisfies

$$\Sigma^* = \mathcal{C}(\Sigma^*) = \sigma_w^2 \mathbb{E}_{z \sim \mathcal{N}(0, \Sigma^*)} [\phi(z) \phi(z)^\top] + \sigma_b^2$$

One solution is perfect correlation,

$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

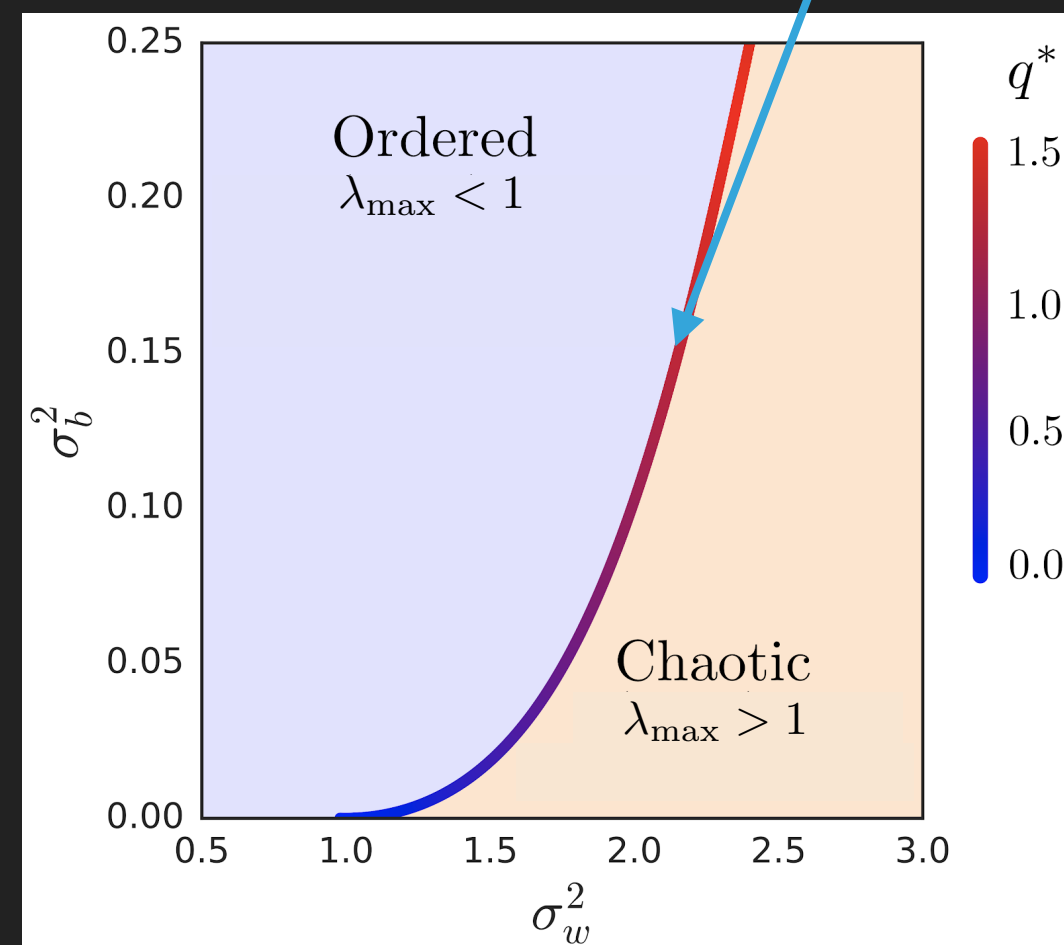
↑ Variance
 ↙ Correlation

Is this fixed point stable?

- Depends on hyperparameters

$$\lambda_{\max} = 1$$

$$\xi = \infty$$



## THE EDGE OF CHAOS

For **deep** signal propagation, initialize on the “edge of chaos”

## THE EDGE OF CHAOS

For **deep** signal propagation, initialize on the “edge of chaos”


- Analyze the fixed point and set  $\lambda_{\max} = 1$



## THE EDGE OF CHAOS

For **deep** signal propagation, initialize on the “edge of chaos”

- Analyze the fixed point and set  $\lambda_{\max} = 1$



$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$


$$\lambda_{\max} \left( \left. \frac{\partial \mathcal{C}(\Sigma)}{\partial \Sigma} \right|_{\Sigma^*} \right) = \chi(\sigma_w, \sigma_b)$$

## THE EDGE OF CHAOS

For **deep** signal propagation, initialize on the “edge of chaos”

- Analyze the fixed point and set  $\lambda_{\max} = 1$


$$\Sigma^* = q^* \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$


$$\lambda_{\max} \left( \left. \frac{\partial \mathcal{C}(\Sigma)}{\partial \Sigma} \right|_{\Sigma^*} \right) = \chi(\sigma_w, \sigma_b)$$

$$\chi(\sigma_w, \sigma_b) = \sigma_w^2 \int dz \frac{e^{-\frac{1}{2}z^2}}{\sqrt{2\pi}} \phi'(\sqrt{q^*}z)^2$$

# BACKPROPAGATED GRADIENTS

Given a loss  $\mathcal{L}$ , back-propagation gives

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1}) \quad \delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} \quad \delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$$

# BACKPROPAGATED GRADIENTS

Given a loss  $\mathcal{L}$ , back-propagation gives

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1}) \quad \delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} \quad \delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$$

Gradients scale like

$$\mathbb{E}[(\delta_1^l)^2] = \mathbb{E}[(\delta_1^{l+1})^2] \sigma_w^2 \mathbb{E}[\phi'(z_1^{l+1})^2]$$

# BACKPROPAGATED GRADIENTS

Given a loss  $\mathcal{L}$ , back-propagation gives

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1}) \quad \delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} \quad \delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$$

Gradients scale like

$$\mathbb{E}[(\delta_1^l)^2] = \mathbb{E}[(\delta_1^{l+1})^2] \underbrace{\sigma_w^2 \mathbb{E}[\phi'(z_1^{l+1})^2]}_{\chi(\sigma_w, \sigma_b)}$$

$$\mathbb{E}[(\delta_1^1)^2] = \mathbb{E}[(\delta_1^L)^2] \chi(\sigma_w, \sigma_b)^L$$

# BACKPROPAGATED GRADIENTS

Given a loss  $\mathcal{L}$ , back-propagation gives

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1}) \quad \delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} \quad \delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$$

Gradients scale like

$$\mathbb{E}[(\delta_1^l)^2] = \mathbb{E}[(\delta_1^{l+1})^2] \underbrace{\sigma_w^2 \mathbb{E}[\phi'(z_1^{l+1})^2]}_{\chi(\sigma_w, \sigma_b)}$$

$$\mathbb{E}[(\delta_1^1)^2] = \mathbb{E}[(\delta_1^L)^2] \chi(\sigma_w, \sigma_b)^L$$

Gradients explode/vanish unless

$$\chi(\sigma_w, \sigma_b) = 1$$

# BACKPROPAGATED GRADIENTS

Given a loss  $\mathcal{L}$ , back-propagation gives

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1}) \quad \delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} \quad \delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$$

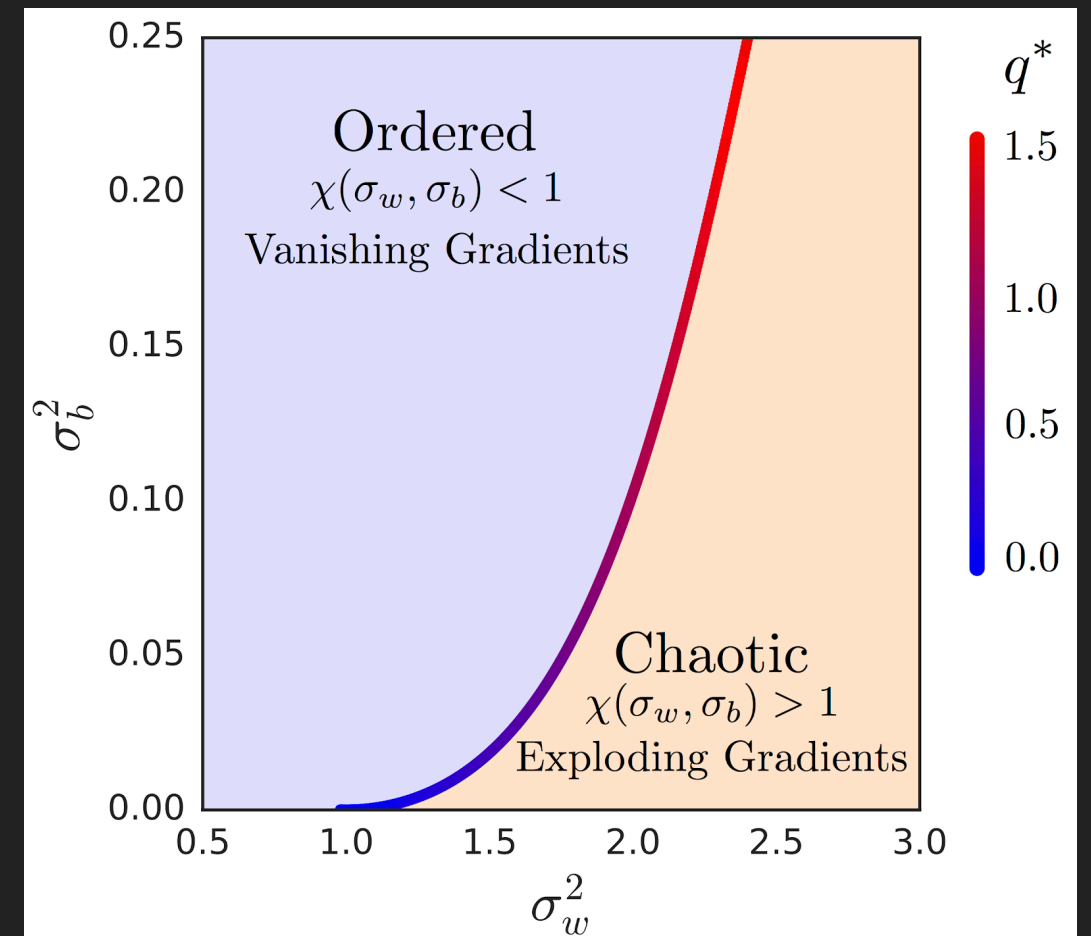
Gradients scale like

$$\mathbb{E}[(\delta_1^l)^2] = \mathbb{E}[(\delta_1^{l+1})^2] \underbrace{\sigma_w^2 \mathbb{E}[\phi'(z_1^{l+1})^2]}_{\chi(\sigma_w, \sigma_b)}$$

$$\mathbb{E}[(\delta_1^1)^2] = \mathbb{E}[(\delta_1^L)^2] \chi(\sigma_w, \sigma_b)^L$$

Gradients explode/vanish unless

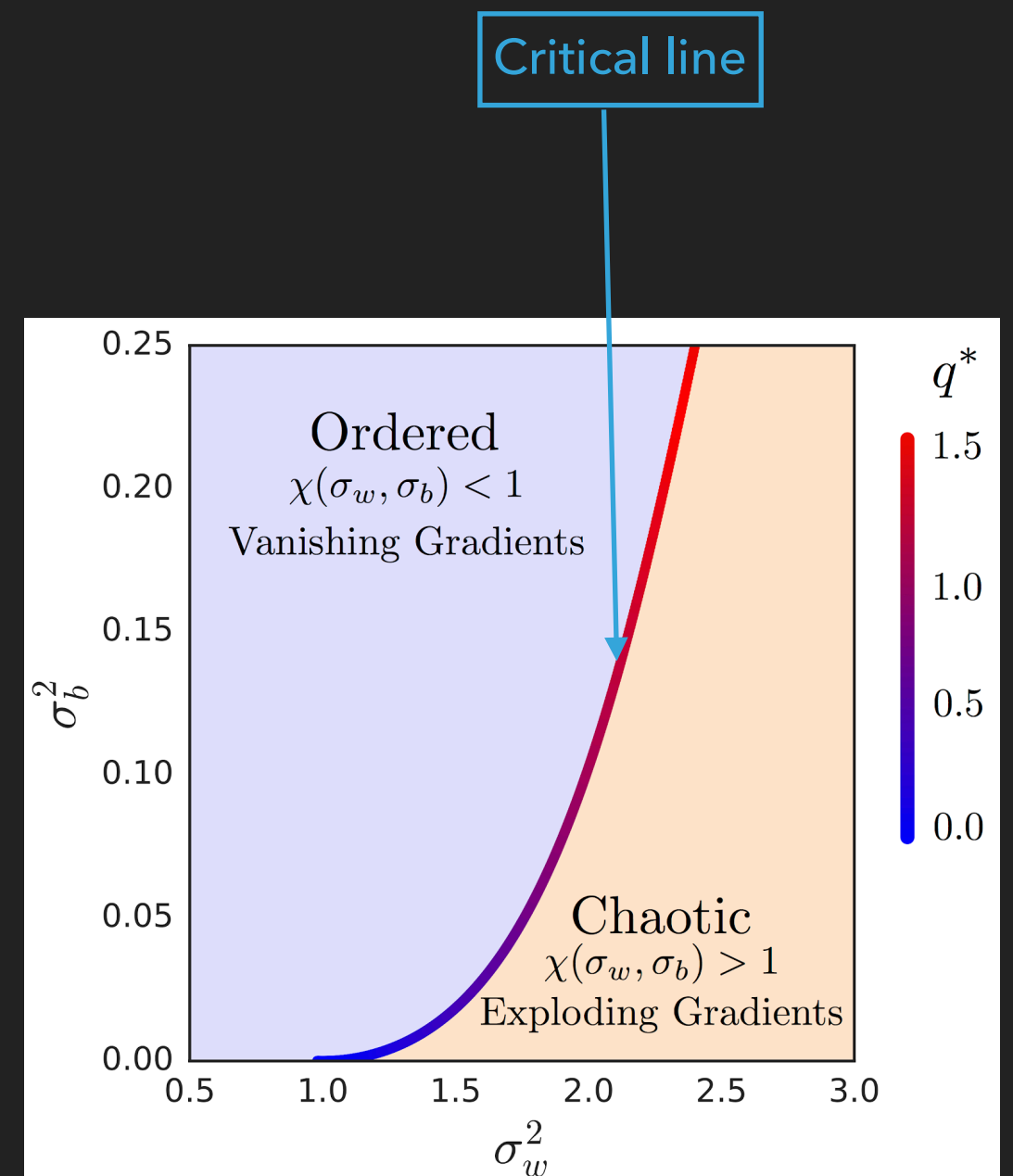
$$\chi(\sigma_w, \sigma_b) = 1$$



# CRITICAL INITIALIZATION

## Critical initialization:

In order for signals to propagate forward and backward through a deep network, the initialization hyperparameters should lie on the critical line

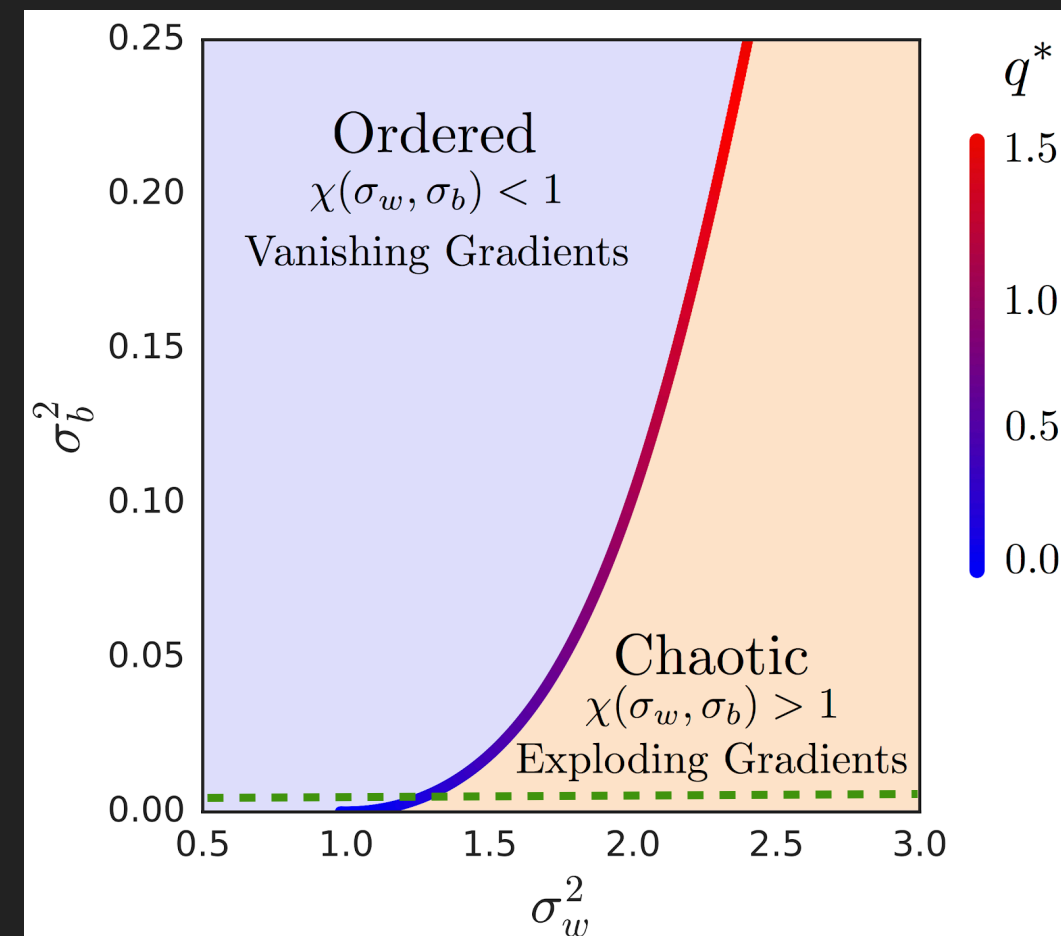
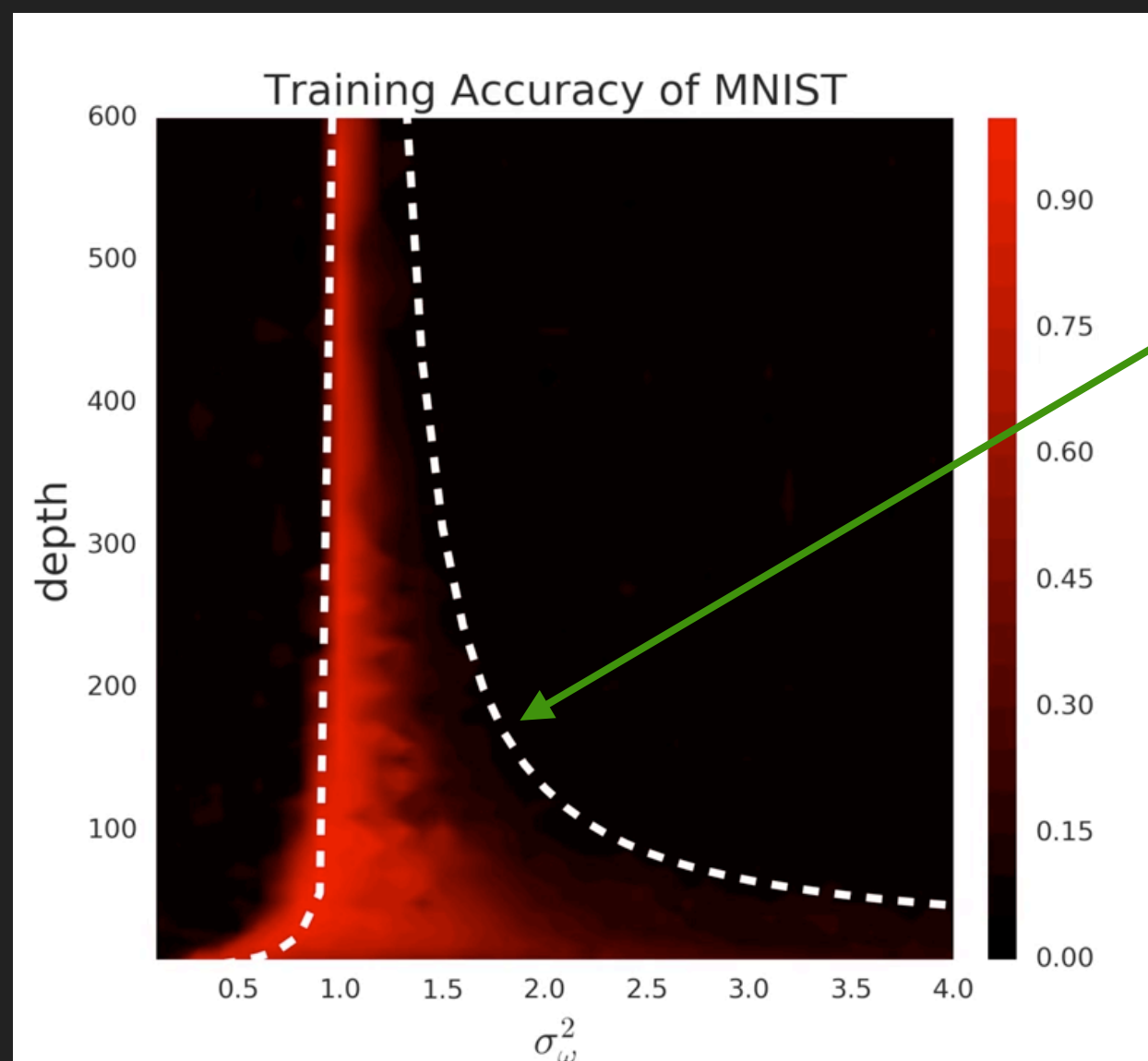




## PREDICTING TRAINABLE DEPTH

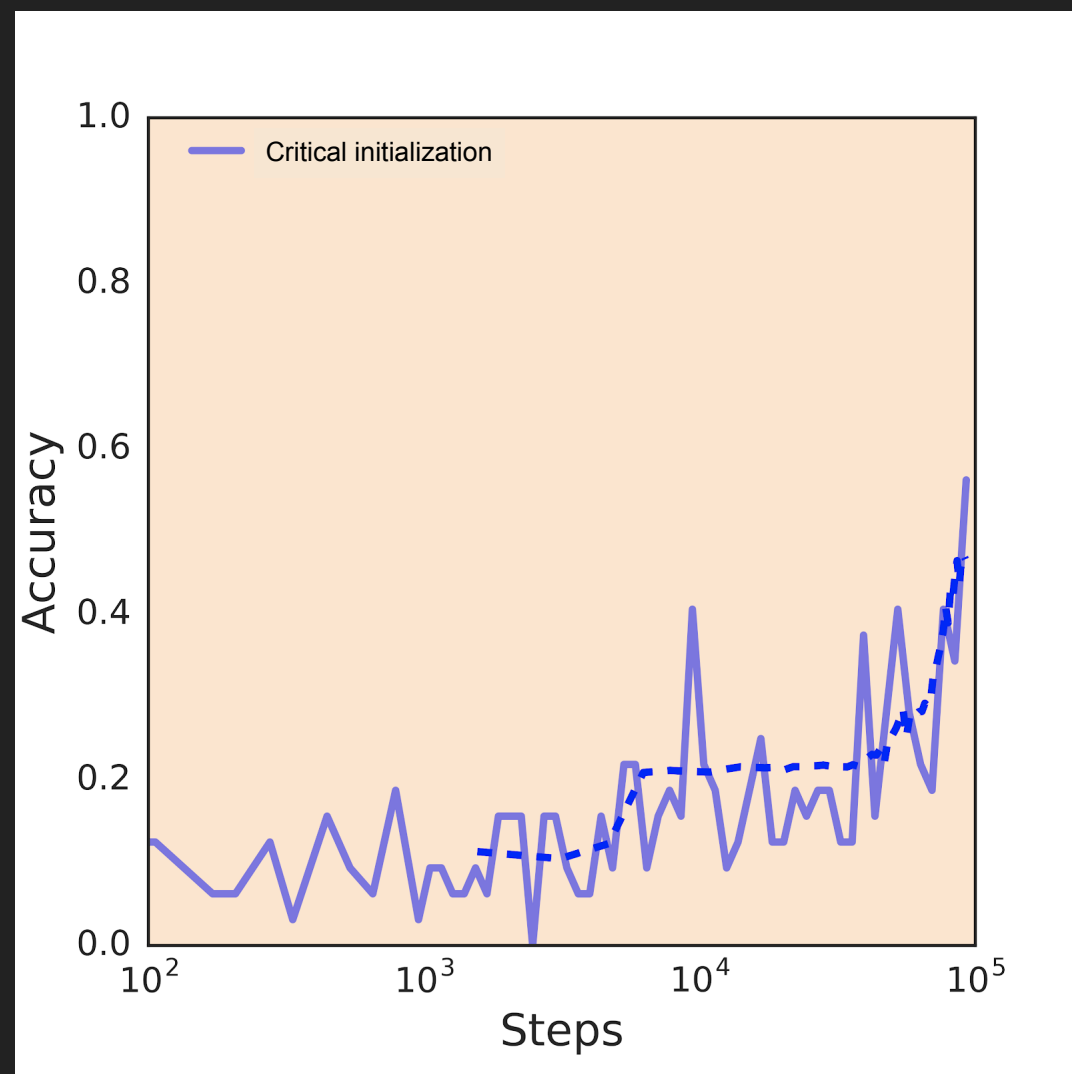
$$\epsilon^l \approx e^{-l/\xi}$$

$$\xi(\sigma_w) = -1 / \log \lambda_{\max}(\sigma_w)$$

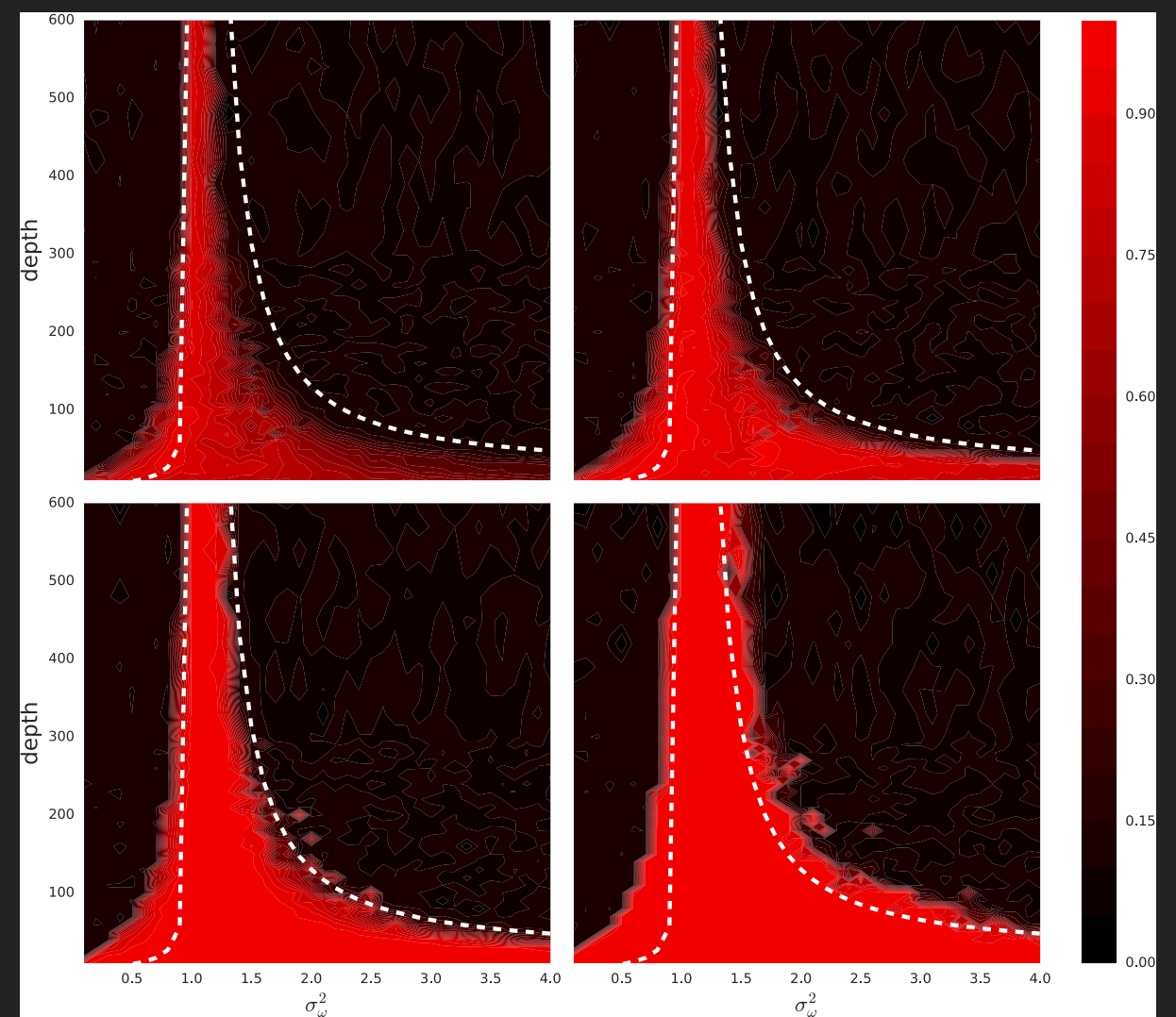


# TRAINABILITY OF VERY DEEP NETWORKS

4000-layer CNN on MNIST



Trainability heat maps




# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion

# DYNAMICAL ISOMETRY

Study the **end-to-end Jacobian**

$$\mathbf{J} = \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^0} = \prod_l \mathbf{D}^l \mathbf{W}^l$$


 Diagonal Matrix

$$D_{ij}^l = \phi'(z_i^l) \delta_{ij}$$

# DYNAMICAL ISOMETRY

Study the **end-to-end Jacobian**

$$\mathbf{J} = \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^0} = \prod_l \mathbf{D}^l \mathbf{W}^l \quad \mathbf{D}_{ij}^l = \phi'(z_i^l) \delta_{ij}$$

 Diagonal Matrix

A few relations that make this interesting

$$\delta^0 = \mathbf{J} \delta^L$$

Gradients

$$f(\mathbf{x} + \delta) \approx f(\mathbf{x}) + \mathbf{J}^T \delta$$

Linear Response

$$\mathbf{G} = \mathbf{J}^T \mathbf{J}$$

Induced Metric

# DYNAMICAL ISOMETRY

Study the **end-to-end Jacobian**

$$\mathbf{J} = \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^0} = \prod_l \mathbf{D}^l \mathbf{W}^l \quad \mathbf{D}_{ij}^l = \phi'(z_i^l) \delta_{ij}$$

↖ Diagonal Matrix

A few relations that make this interesting

$$\delta^0 = \mathbf{J} \delta^L \quad f(\mathbf{x} + \delta) \approx f(\mathbf{x}) + \mathbf{J}^T \delta \quad G = \mathbf{J}^T \mathbf{J}$$

Gradients                      Linear Response                      Induced Metric

We have worked out behavior of gradients **on average**:

$$\text{Criticality} \quad \Leftrightarrow \quad \mathbb{E}[\text{tr}(\mathbf{J}^T \mathbf{J})] = \chi(\sigma_w, \sigma_b)^L = 1$$

# DYNAMICAL ISOMETRY

Study the **end-to-end Jacobian**

$$\mathbf{J} = \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^0} = \prod_l \mathbf{D}^l \mathbf{W}^l \quad \mathbf{D}_{ij}^l = \phi'(z_i^l) \delta_{ij}$$

↖ Diagonal Matrix

A few relations that make this interesting

$$\delta^0 = \mathbf{J} \delta^L \quad f(\mathbf{x} + \delta) \approx f(\mathbf{x}) + \mathbf{J}^T \delta \quad G = \mathbf{J}^T \mathbf{J}$$

Gradients                      Linear Response                      Induced Metric

We have worked out behavior of gradients **on average**:

$$\text{Criticality} \quad \Leftrightarrow \quad \mathbb{E}[\text{tr}(\mathbf{J}^T \mathbf{J})] = \chi(\sigma_w, \sigma_b)^L = 1$$

But what is a good prior for the whole spectrum?

# DYNAMICAL ISOMETRY

Study the **end-to-end Jacobian**

$$\mathbf{J} = \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^0} = \prod_l \mathbf{D}^l \mathbf{W}^l \quad \mathbf{D}_{ij}^l = \phi'(z_i^l) \delta_{ij}$$

↖ Diagonal Matrix

A few relations that make this interesting

$$\delta^0 = \mathbf{J} \delta^L \quad f(\mathbf{x} + \delta) \approx f(\mathbf{x}) + \mathbf{J}^T \delta \quad G = \mathbf{J}^T \mathbf{J}$$

Gradients                      Linear Response                      Induced Metric

We have worked out behavior of gradients **on average**:

$$\text{Criticality} \quad \Leftrightarrow \quad \mathbb{E}[\text{tr}(\mathbf{J}^T \mathbf{J})] = \chi(\sigma_w, \sigma_b)^L = 1$$

But what is a good prior for the whole spectrum?

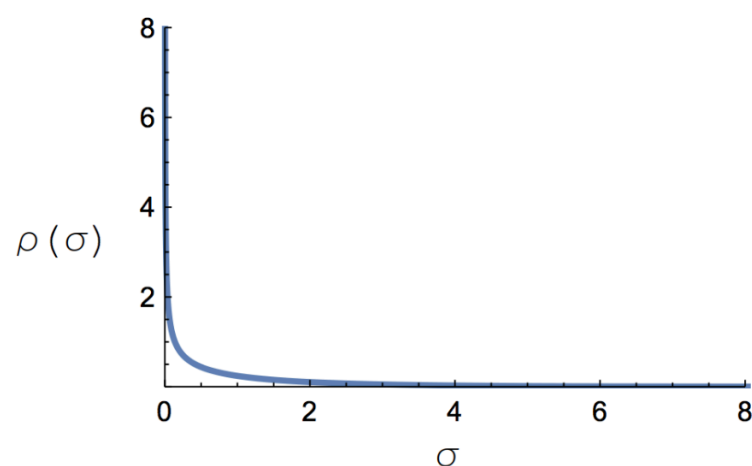
- **Isometry**: all singular values  $\approx 1$



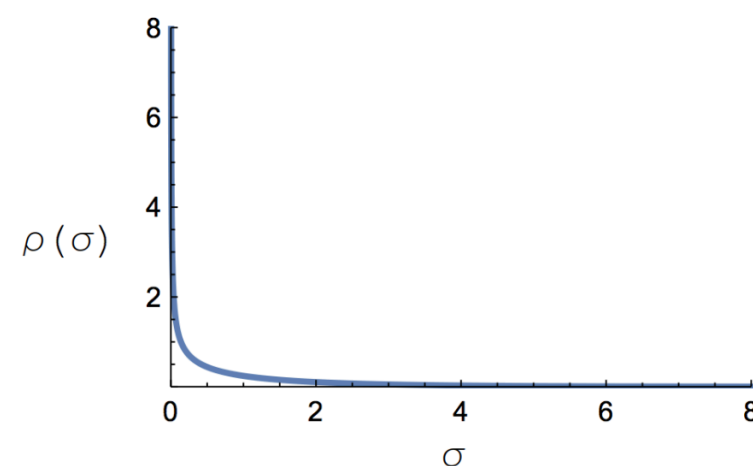
# COMPUTING THE SPECTRUM

Using tools from random matrix theory (free probability), can compute spectrum analytically:

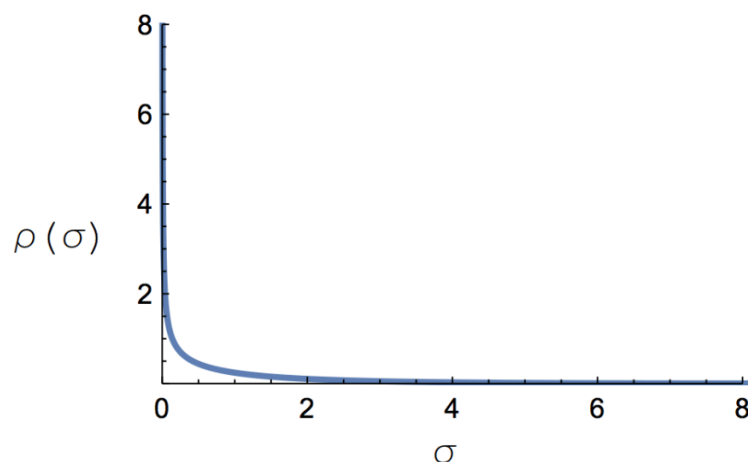
Gaussian W, any f



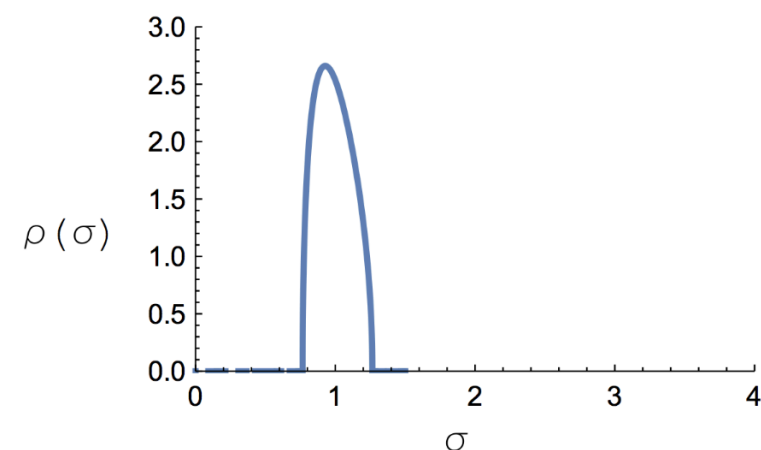
Orthogonal W, ReLU



Orthogonal W, tanh,  $\sigma_w \gg 1$



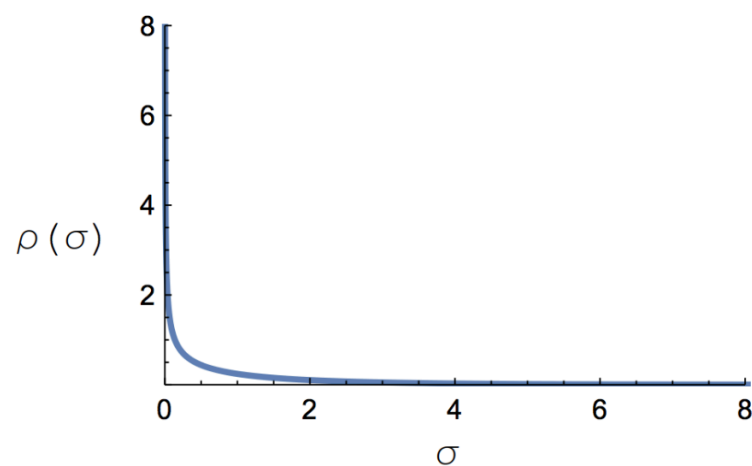
Orthogonal W, tanh,  $\sigma_w \sim 1+1/L$



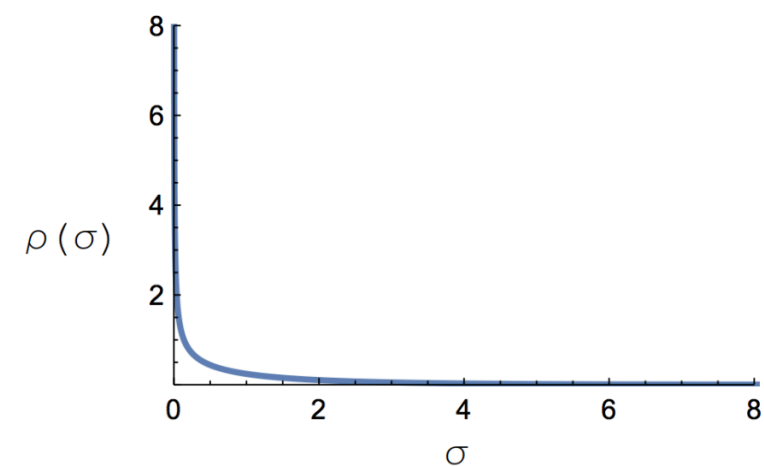
# COMPUTING THE SPECTRUM

Using tools from random matrix theory (free probability), can compute spectrum analytically:

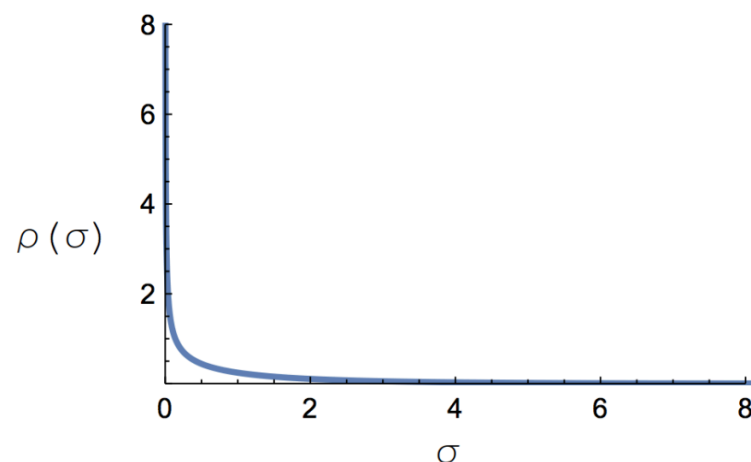
Gaussian W, any f



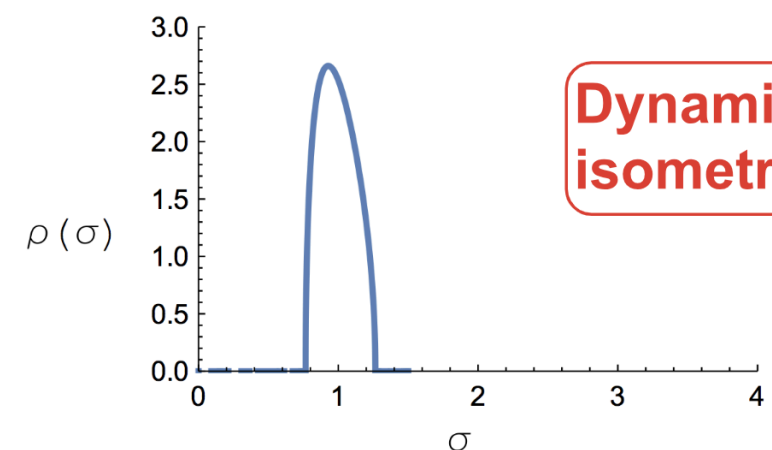
Orthogonal W, ReLU



Orthogonal W, tanh,  $\sigma_w \gg 1$



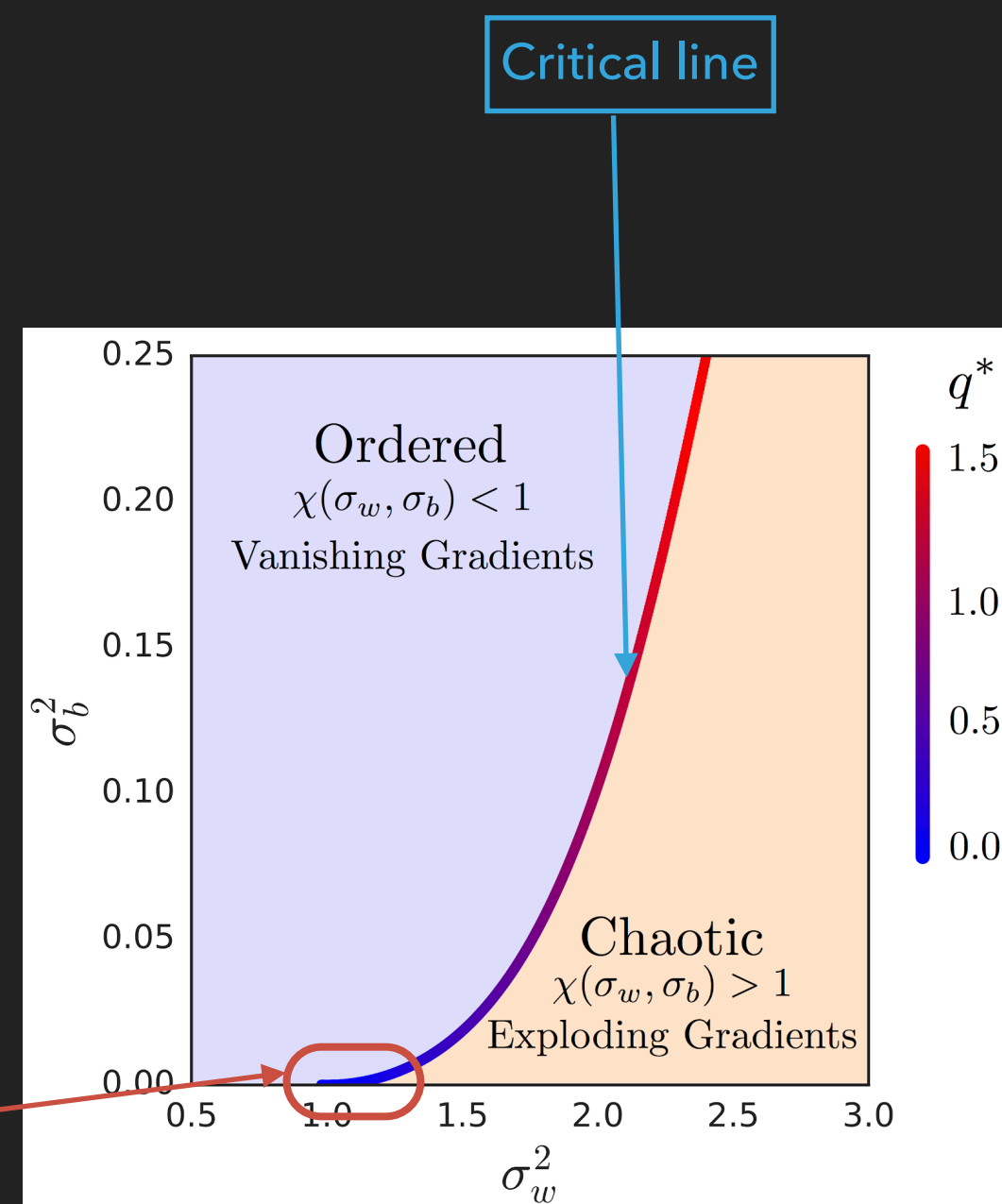
Orthogonal W, tanh,  $\sigma_w \sim 1+1/L$



# DYNAMICAL ISOMETRY

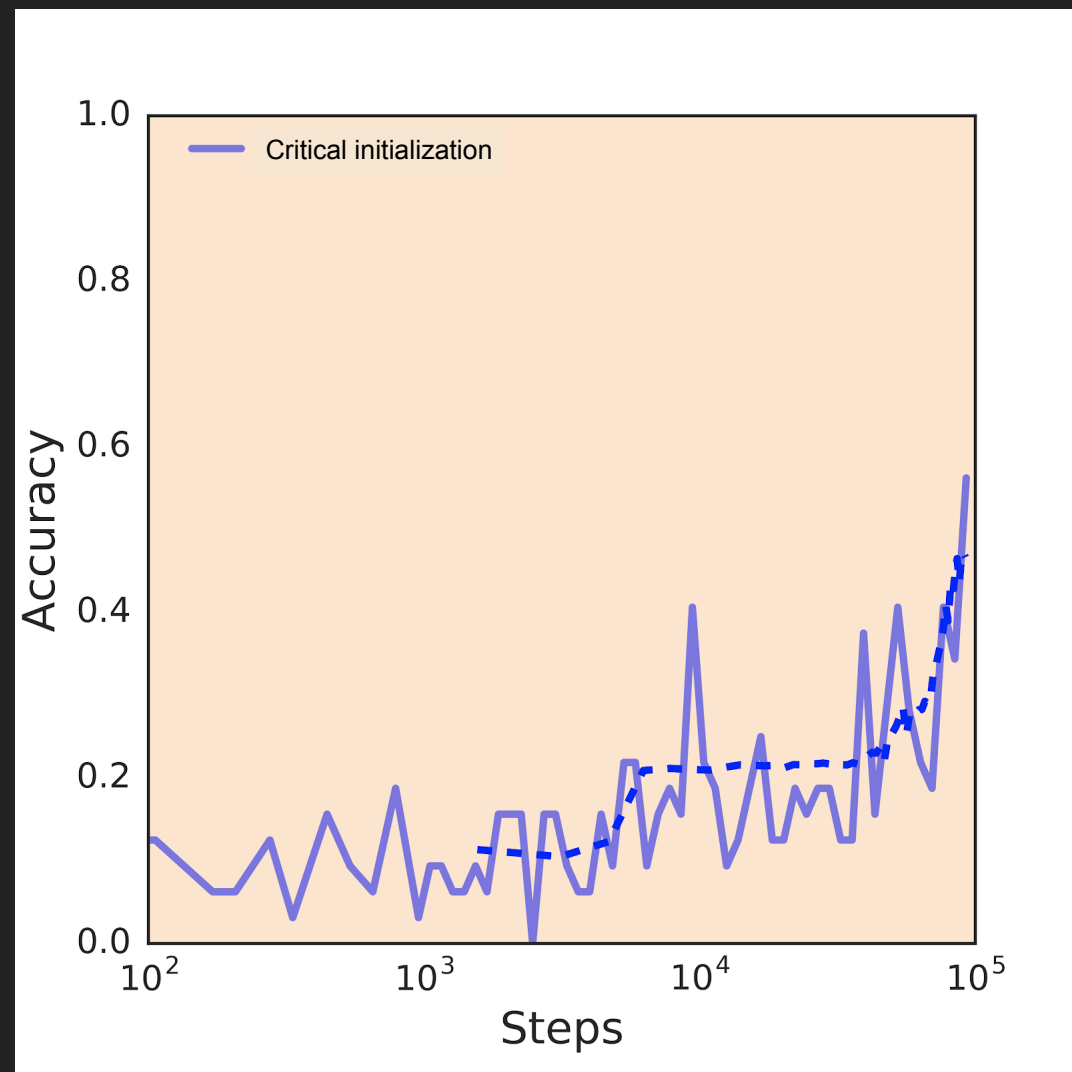
Not every point on the critical line is equally favorable for gradient propagation. For activation functions that are linear near the origin, **dynamical isometry** (i.e. well-conditioned Jacobians) can be achieved with small bias variance.

Dynamical isometry



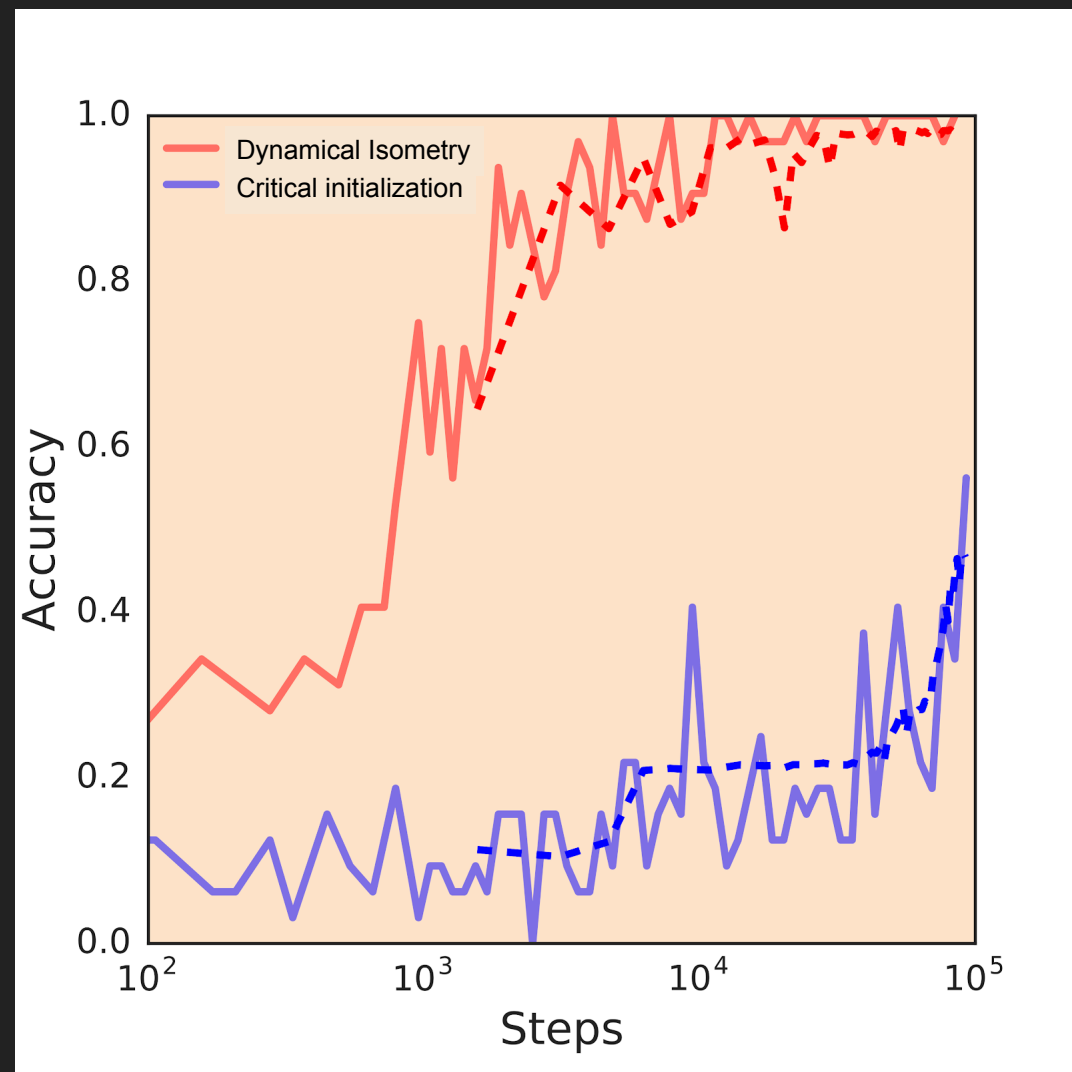
# THE BENEFITS OF A BETTER PRIOR

4000-layer CNN on MNIST



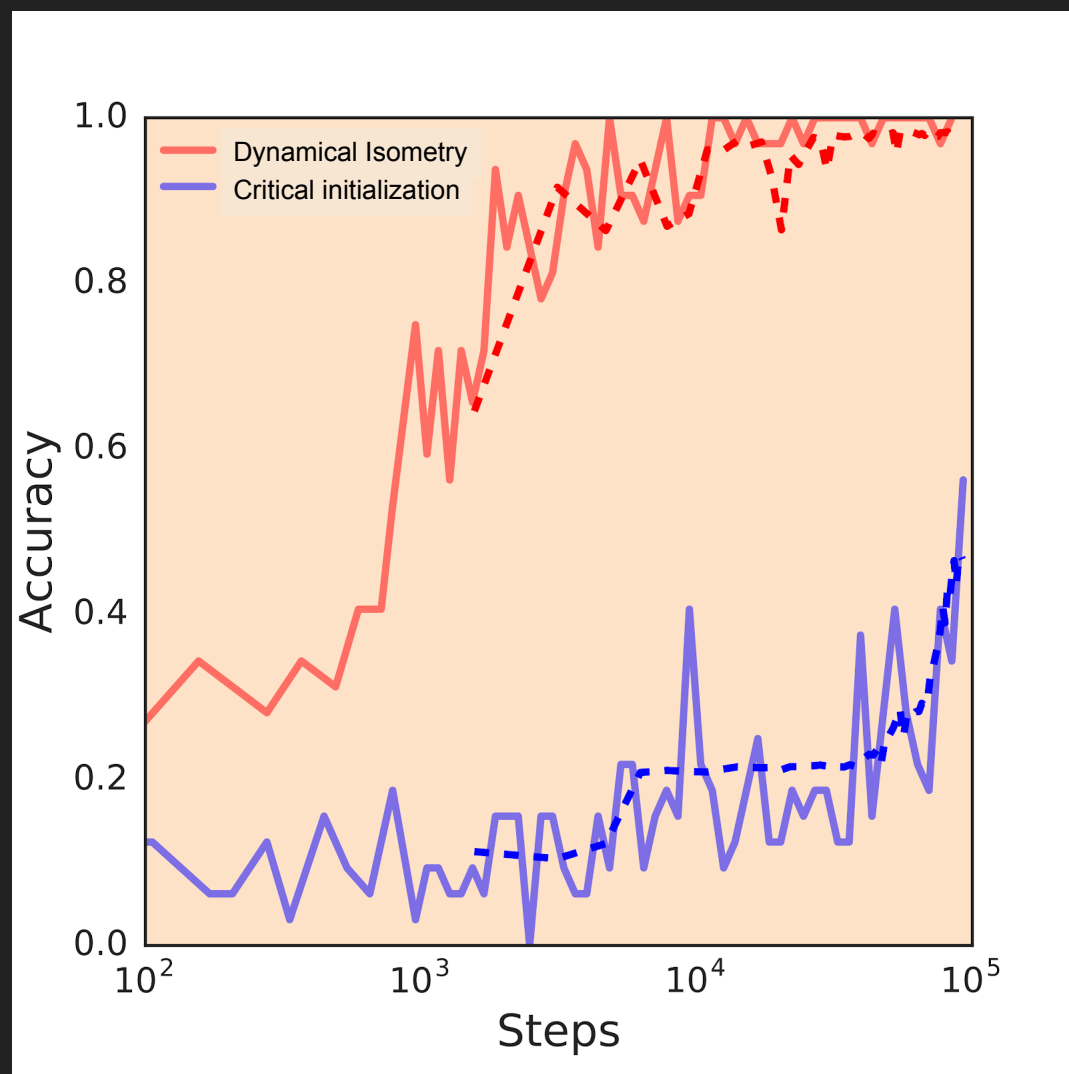
# THE BENEFITS OF A BETTER PRIOR

4000-layer CNN on MNIST

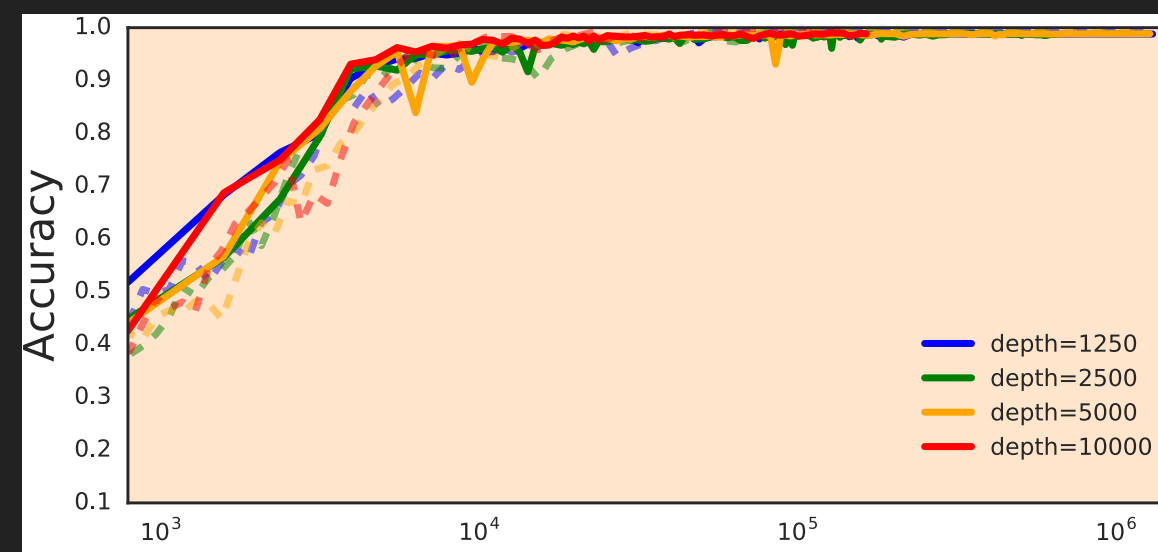


# THE BENEFITS OF A BETTER PRIOR

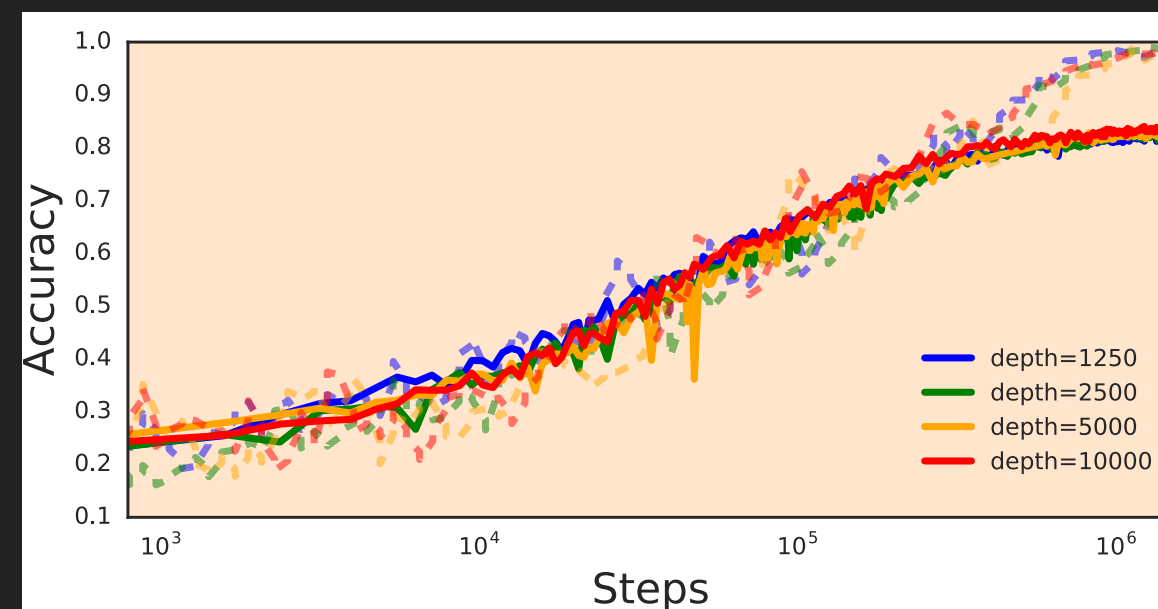
## 4000-layer CNN on MNIST



## MNIST



## CIFAR-10




# OUTLINE

1. Motivation
2. Functional priors
3. Signal propagation
4. Dynamical isometry
5. Functional posteriors
6. Conclusion

# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Consider a FC neural network,  $f(x; \theta(t))$

  
Gradient Descent Time



# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Consider a FC neural network,  $f(x; \theta(t))$



Gradient Descent Time

As the width of the network grows, parameters move less during gradient descent

# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Consider a FC neural network,  $f(x; \theta(t))$

Gradient Descent Time

As the width of the network grows, parameters move less during gradient descent

Motivates a linear approximation,

$$f(x; \theta(t)) \approx f(x; \theta(0)) + \sum_{\alpha} \frac{\partial f(x; \theta(0))}{\partial \theta_{\alpha}(0)} (\theta(t) - \theta(0)) + \mathcal{O}((\theta(t) - \theta(0))^2)$$

Function at Initialization

Jacobian at Initialization

This becomes exact as  $N \rightarrow \infty$

# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Consider a FC neural network,  $f(x; \theta(t))$

Gradient Descent Time

As the width of the network grows, parameters move less during gradient descent

Motivates a linear approximation,

$$f_t(x) \approx f_0(x) + J_0(x)\omega(t) \qquad \omega(t) = \theta(t) - \theta(0)$$

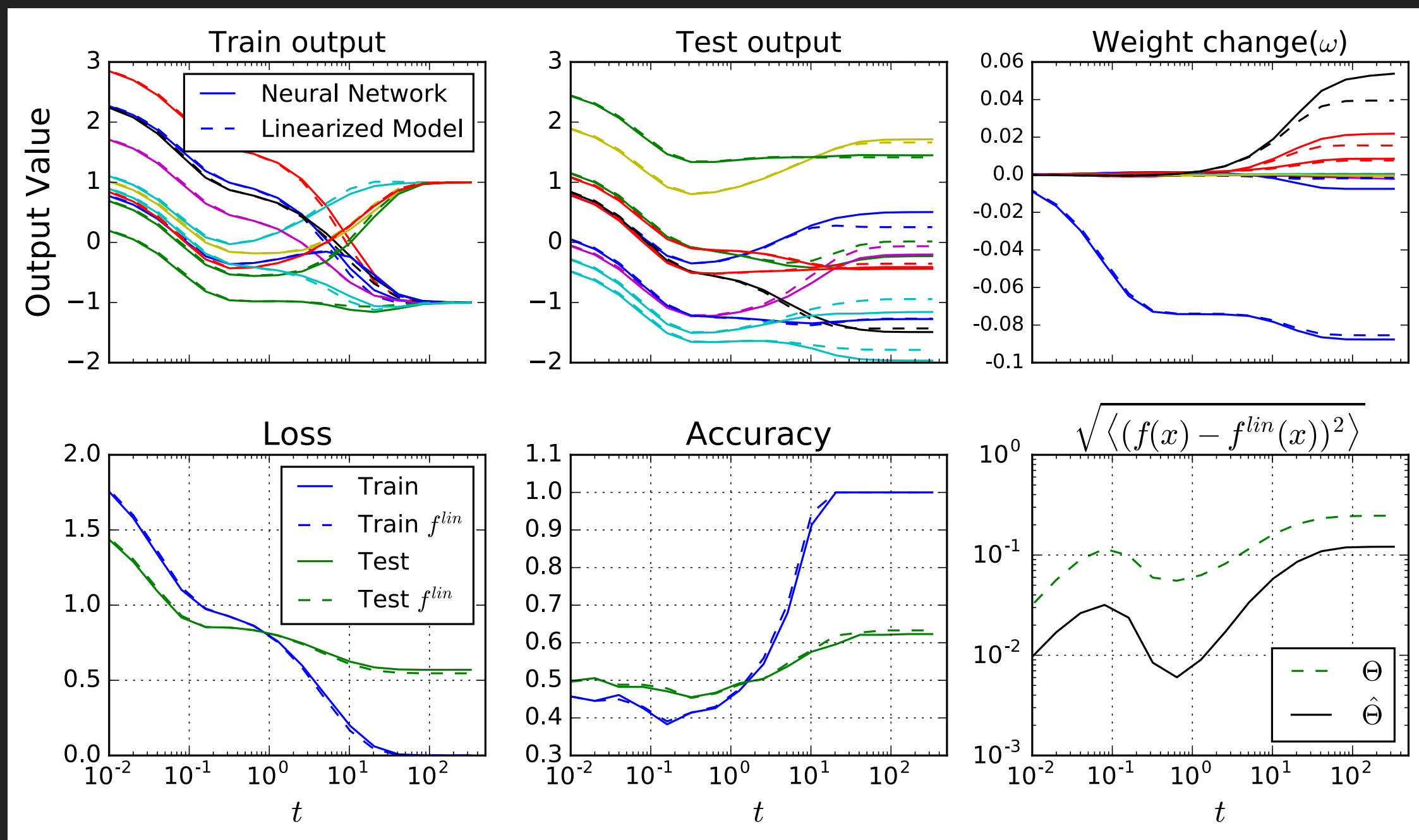
Function at Initialization

Jacobian at Initialization

This becomes exact as  $N \rightarrow \infty$

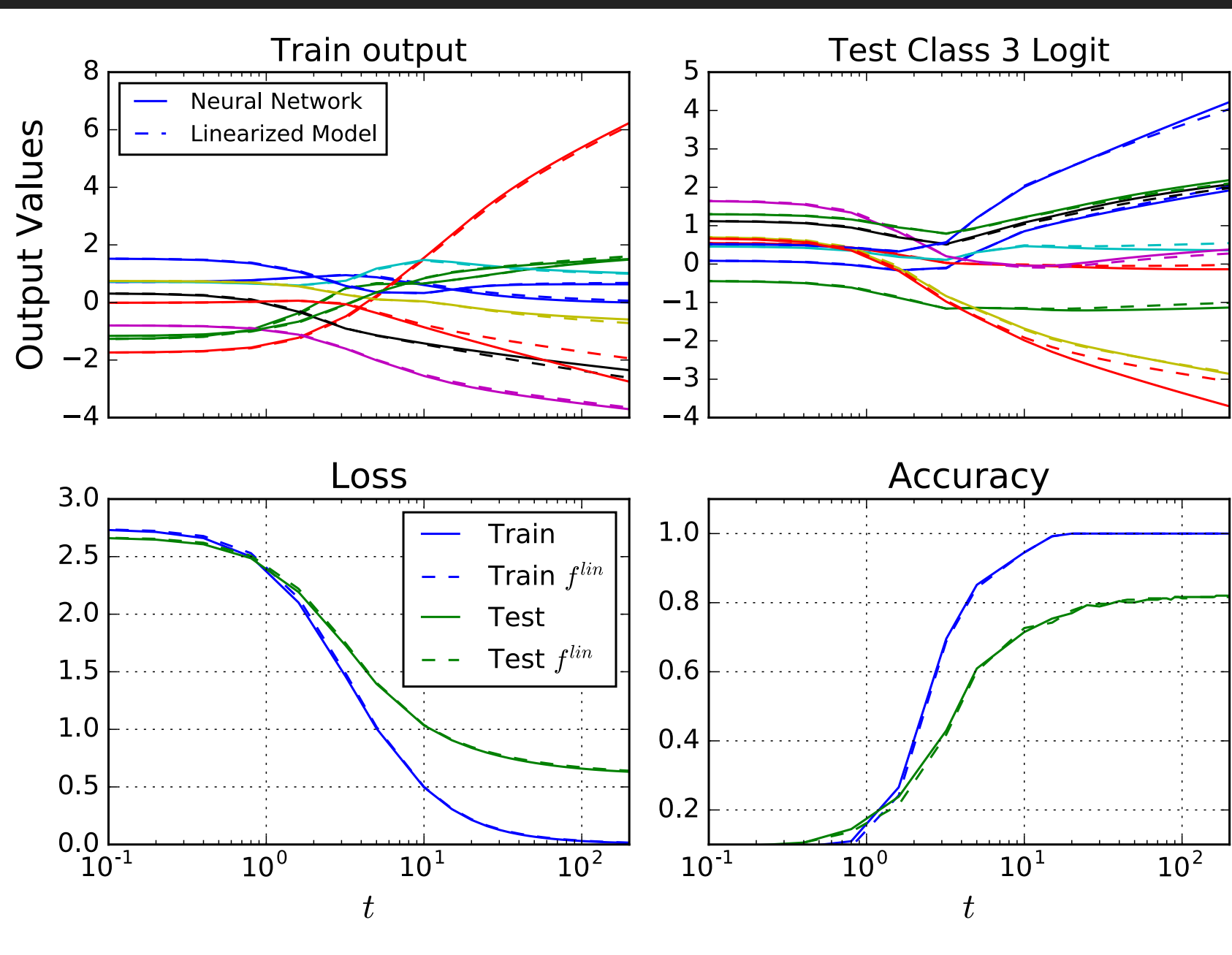
# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Fully Connected, N=2048, Single Output, MSE Loss, Gradient Descent



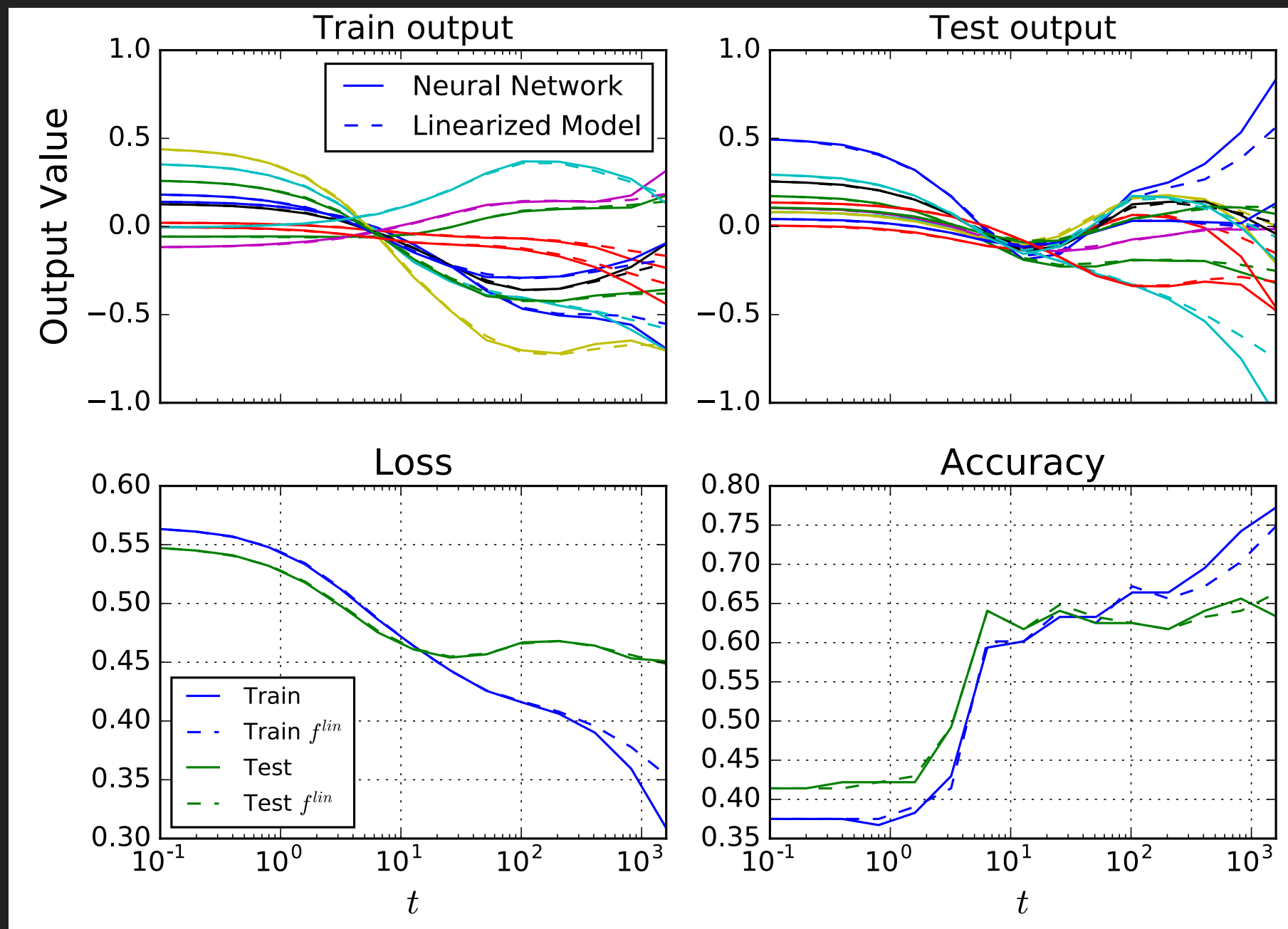
# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Fully Connected, N=1024, 10-Class, Cross Entropy Loss, Momentum



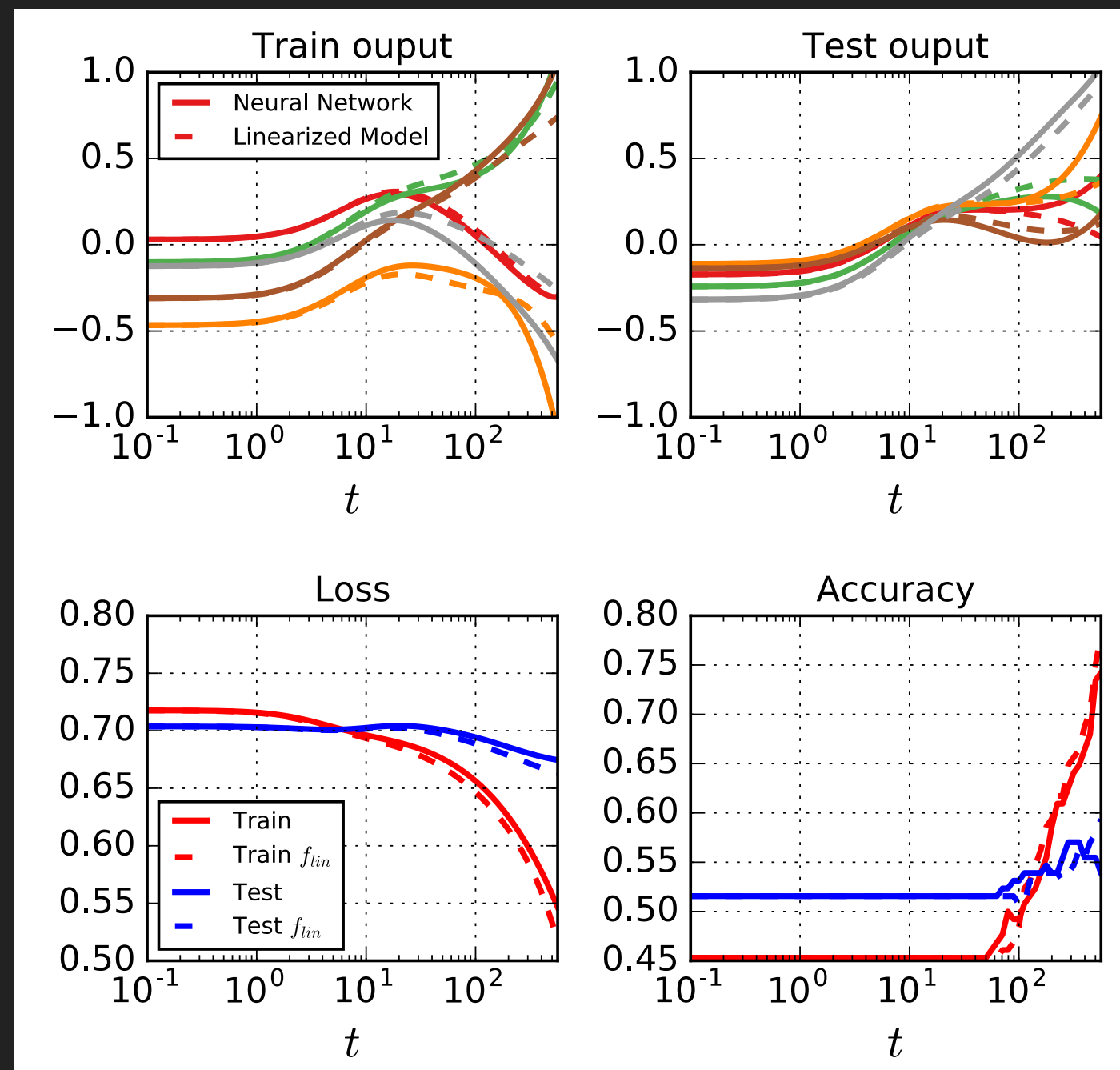
# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

CNN, C=256, 2-Class, MSE Loss, GD



# WIDE, DEEP, NETWORKS EVOLVE AS LINEAR MODELS

Wide Resnet (10-layers), C=1024, 2-Class, Cross Entropy Loss, Momentum



# IMPLICATIONS FOR THE POSTERIOR

For MSE Loss,

$$\partial_t f_t(X) = -\Theta(X, X)(f_t(X) - Y) \quad \Theta(X, Y) = \frac{1}{M} J_0(X) J_0(Y)^T$$


↑  
Neural Tangent Kernel



# IMPLICATIONS FOR THE POSTERIOR

For MSE Loss,


$$\partial_t f_t(X) = -\Theta(X, X)(f_t(X) - Y) \quad \Theta(X, Y) = \frac{1}{M} J_0(X) J_0(Y)^T$$


 Neural Tangent Kernel

This allows us to compute the “posterior” after  $t$  steps of GD,

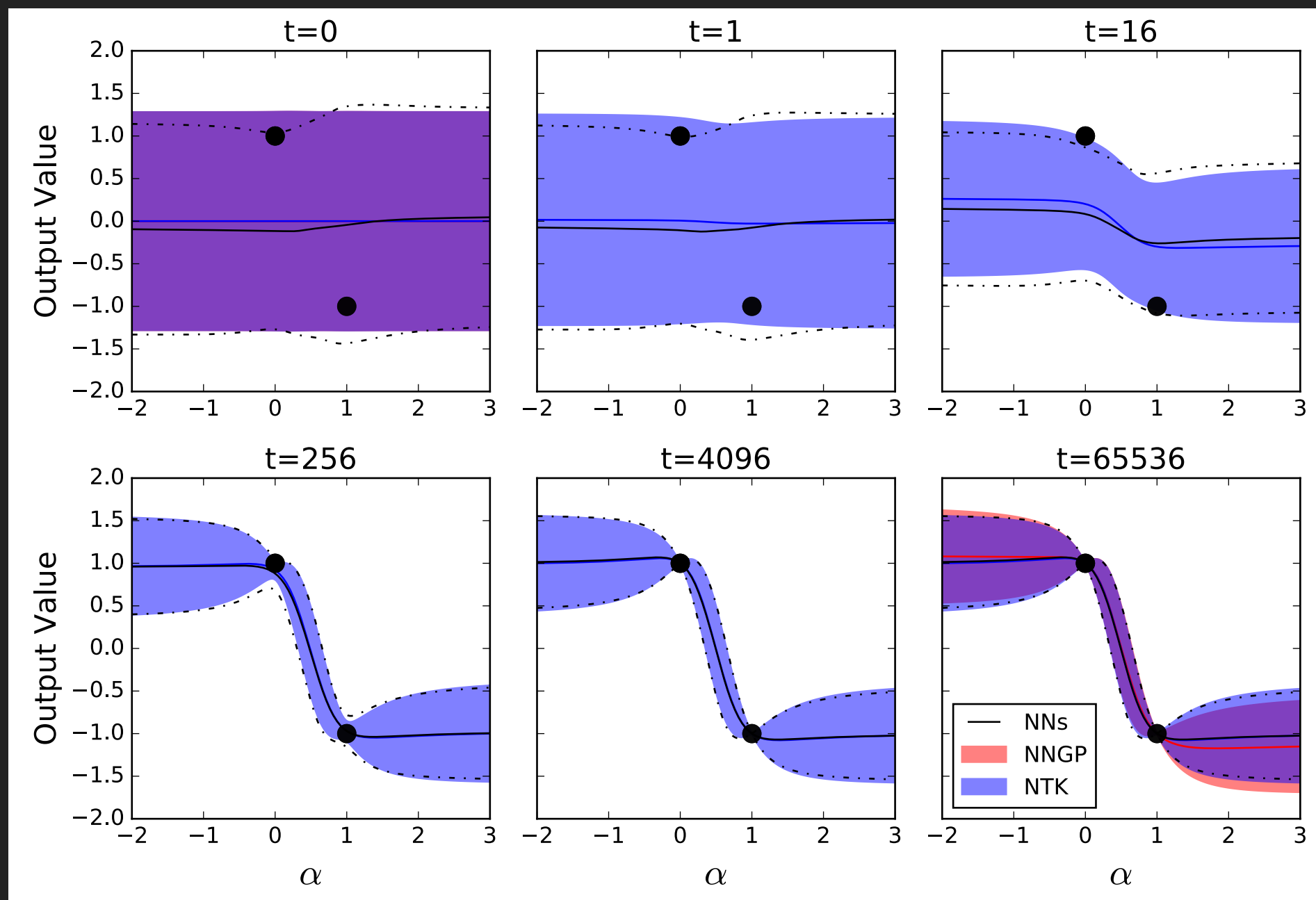
$$\mu(x) = \Theta(x, X) \Theta^{-1} (I - e^{-\eta \Theta t}) Y$$

$$\begin{aligned} \Sigma(x) = & K(x, x) - 2\Theta(x, X) \Theta^{-1} (I - e^{-\eta \Theta t}) K(x, X)^T \\ & + \Theta(x, X) \Theta^{-1} (I - e^{-\eta \Theta t}) K \Theta^{-1} (I - e^{-\eta \Theta t}) \Theta(x, X)^T \end{aligned}$$


 NNGP Kernel

# IMPLICATIONS FOR THE POSTERIOR

FC Network,  $N=8192$ , MNIST, MSE Loss



# CONCLUSIONS

Overparameterized models are simple!

The prior over functions can be computed analytically

Properties of the prior are intimately related to trainability

Wide neural networks are almost linear models

Overall, a powerful framework is emerging for theoretically analyzing overparameterized neural networks