

# Effective Theory of Deep Neural Networks

Sho Yaida

# Effective Theory of Deep Neural Networks

Dan Roberts

Sho Yaida

Boris Hanin

[\[arXiv:2106.10165\]](https://arxiv.org/abs/2106.10165)

# Effective Theory of Deep Neural Networks

Dan Roberts

Sho Yaida

Boris Hanin

[[arXiv:2106.10165](https://arxiv.org/abs/2106.10165)~470 pages]

# Agenda

1. Overview
2. Neural Networks at Infinite Width
3. Neural Networks at Finite Width
4. The Principles

# 1. Overview

Deep learning is powerful

Deep learning is powerful

**[put an  
impressive chart]**

**[put an  
impressive picture]**

**[put an  
impressive text]**

Deep learning is powerful & interesting

**[put an  
impressive chart]**

**[put an  
impressive picture]**

**[put an  
impressive text]**



# Machine Learning in a Nutshell

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \theta_{\text{init}} \in p(\theta_{\text{init}})$$

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \theta_{\text{init}} \in p(\theta_{\text{init}})$$

- Train the model, e.g. by gradient descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \right|_{\theta = \theta(t)}$$

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \theta_{\text{init}} \in p(\theta_{\text{init}})$$

- Train the model, e.g. by gradient descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \right|_{\theta = \theta(t)}$$

- Use the trained model to make predictions

$$f_{\text{trained}}(x) = f(x; \theta_{\text{trained}})$$

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \theta_{\text{init}} \in p(\theta_{\text{init}})$$

- Train the model, e.g. by gradient descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \right|_{\theta = \theta(t)}$$

- Use the trained model to make predictions

$$f_{\text{trained}}(x) = f(x; \theta_{\text{trained}})$$

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \underline{\theta_{\text{init}} \in p(\theta_{\text{init}})}$$

- Train the model, e.g. by gradient descent

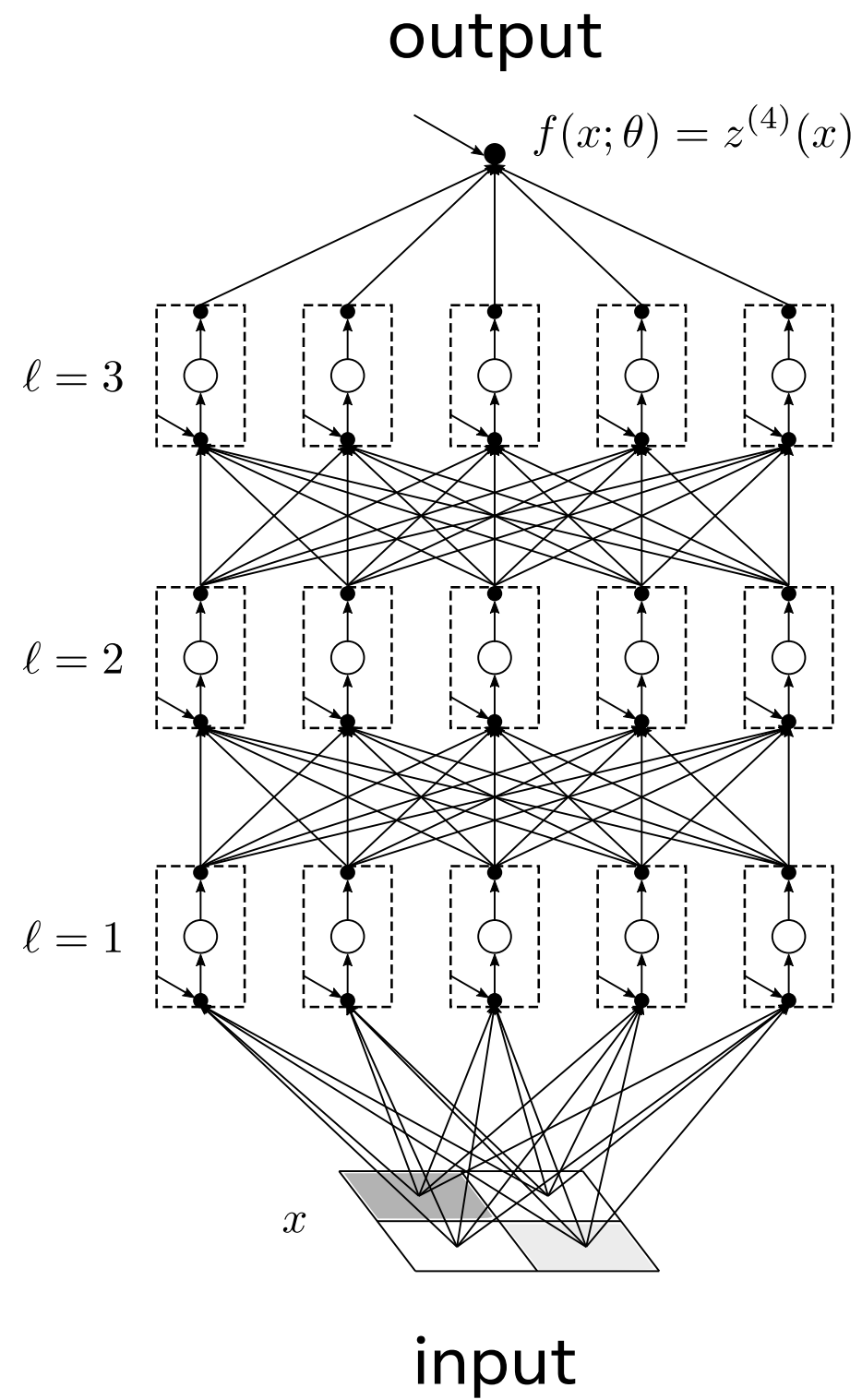
$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \right|_{\theta = \theta(t)}$$

- Use the trained model to make predictions

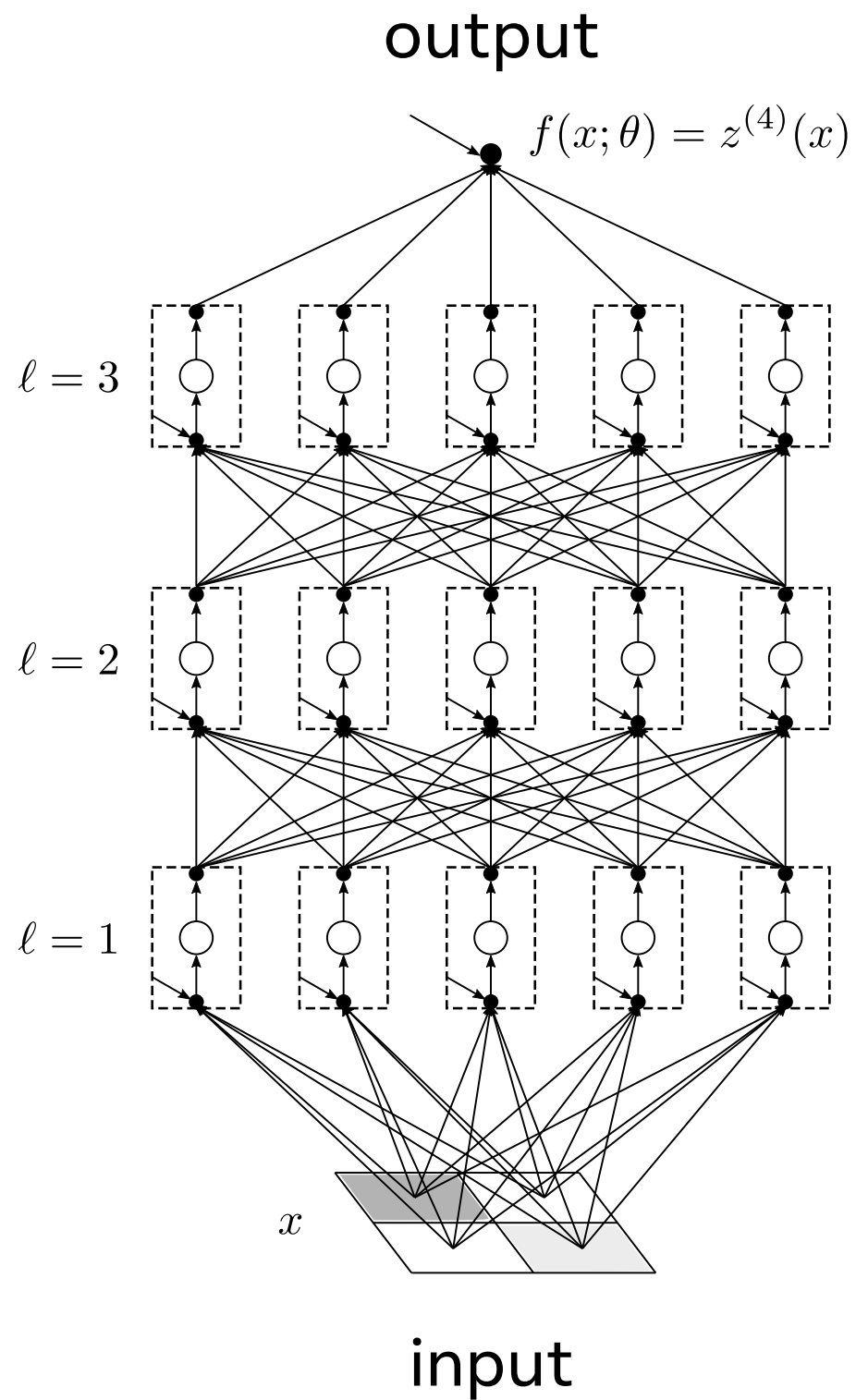
$$p(f_{\text{trained}})$$

mean, variance, etc.

# Neural Networks



# Neural Networks



- Function:

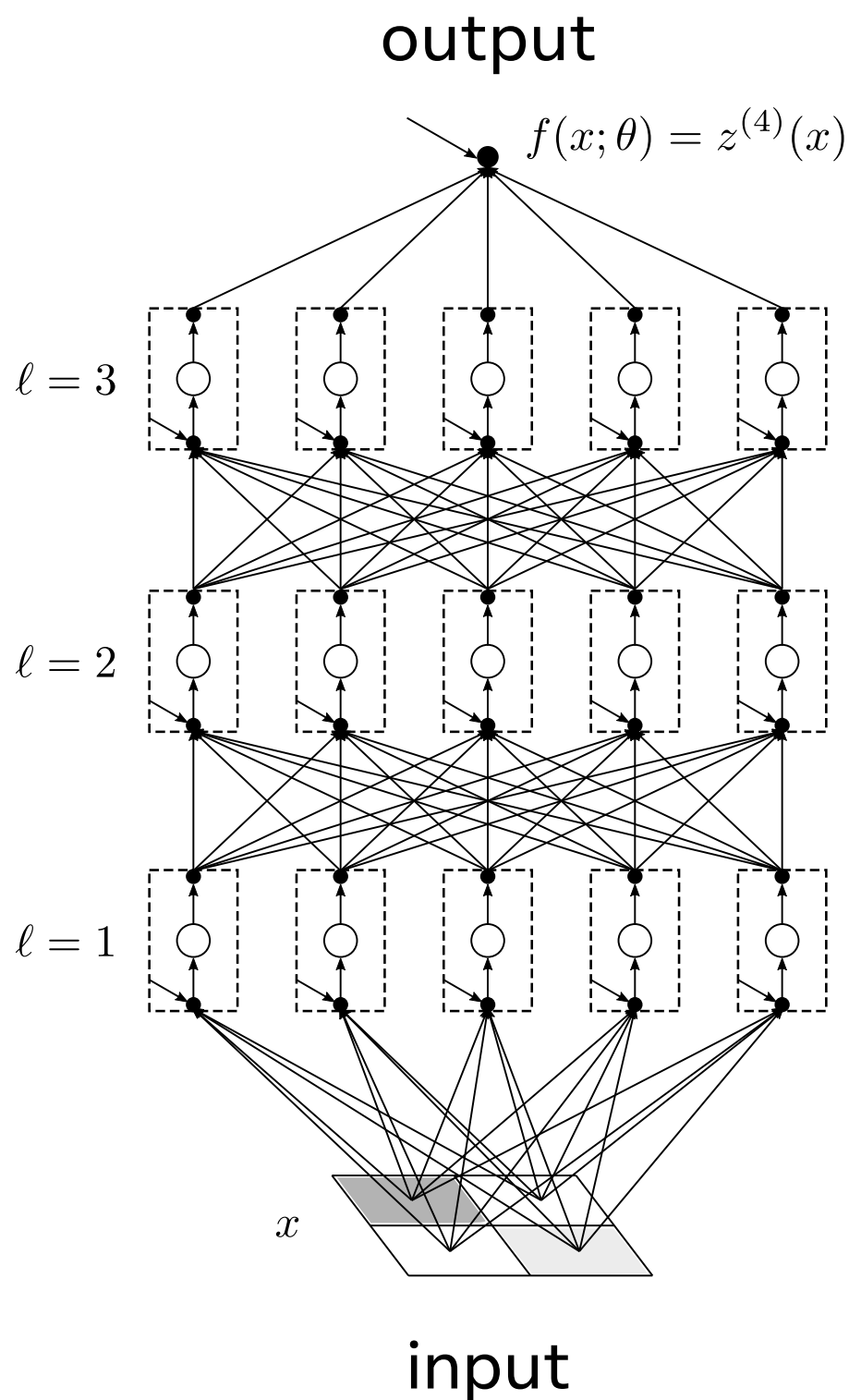
$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)}(x)) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$f(x; \theta) = z^{(L)}(x)$$



# Neural Networks



- Function:

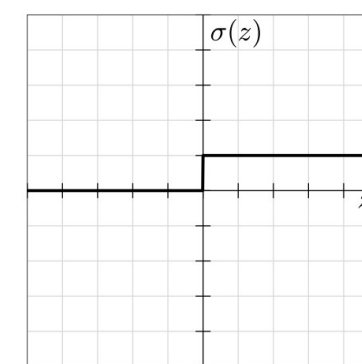
$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)}(x)) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

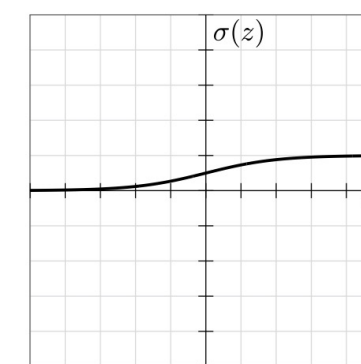
$$f(x; \theta) = z^{(L)}(x)$$

activation function

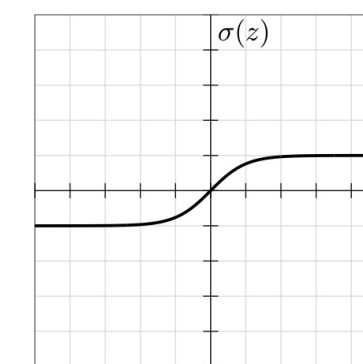
$$\sigma(z)$$



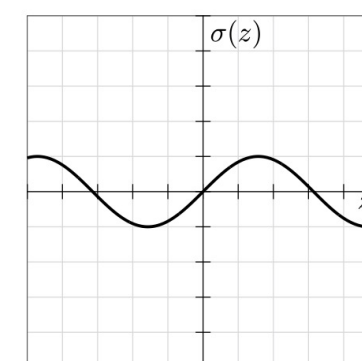
perceptron



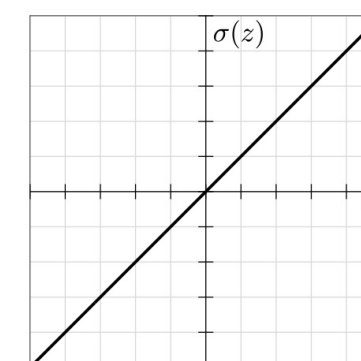
sigmoid



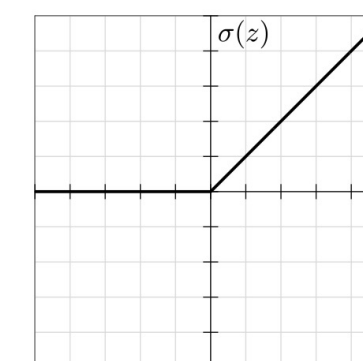
tanh



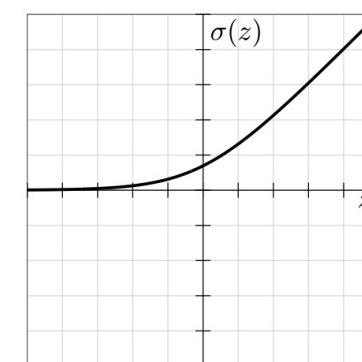
sin



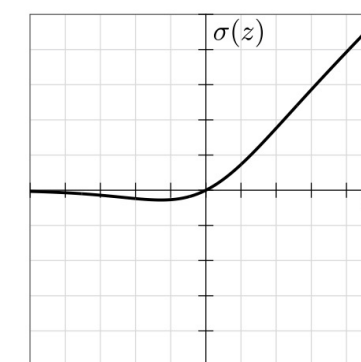
linear



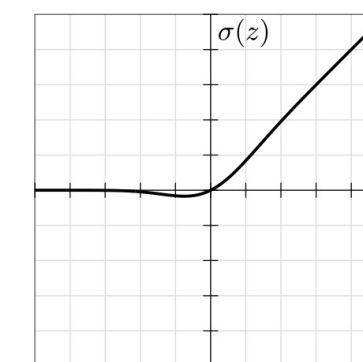
ReLU



softplus

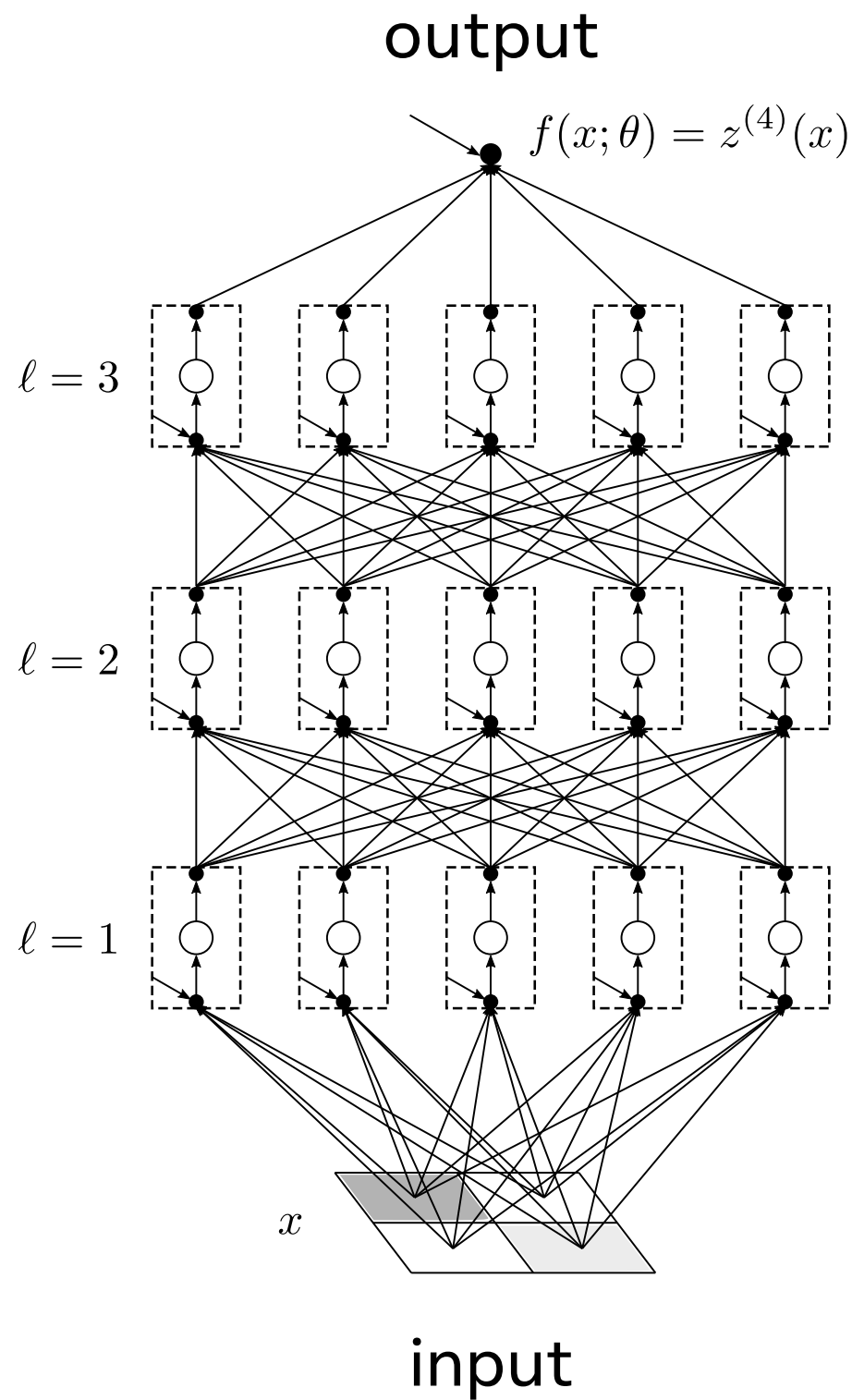


SWISH



GELU

# Neural Networks



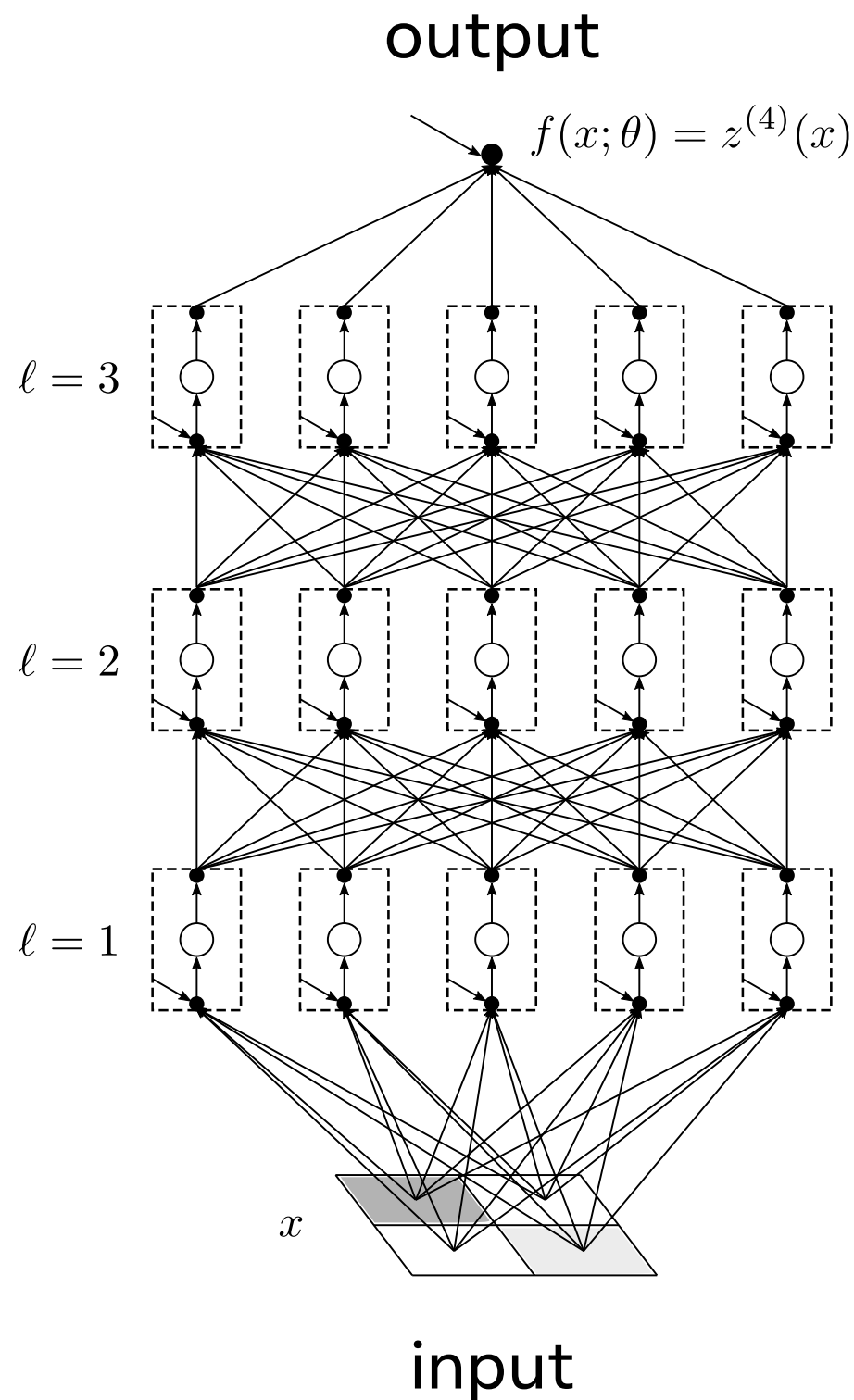
- Function:

$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(z_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$f(x; \theta) = z^{(L)}(x)$$

# Neural Networks



- Function:

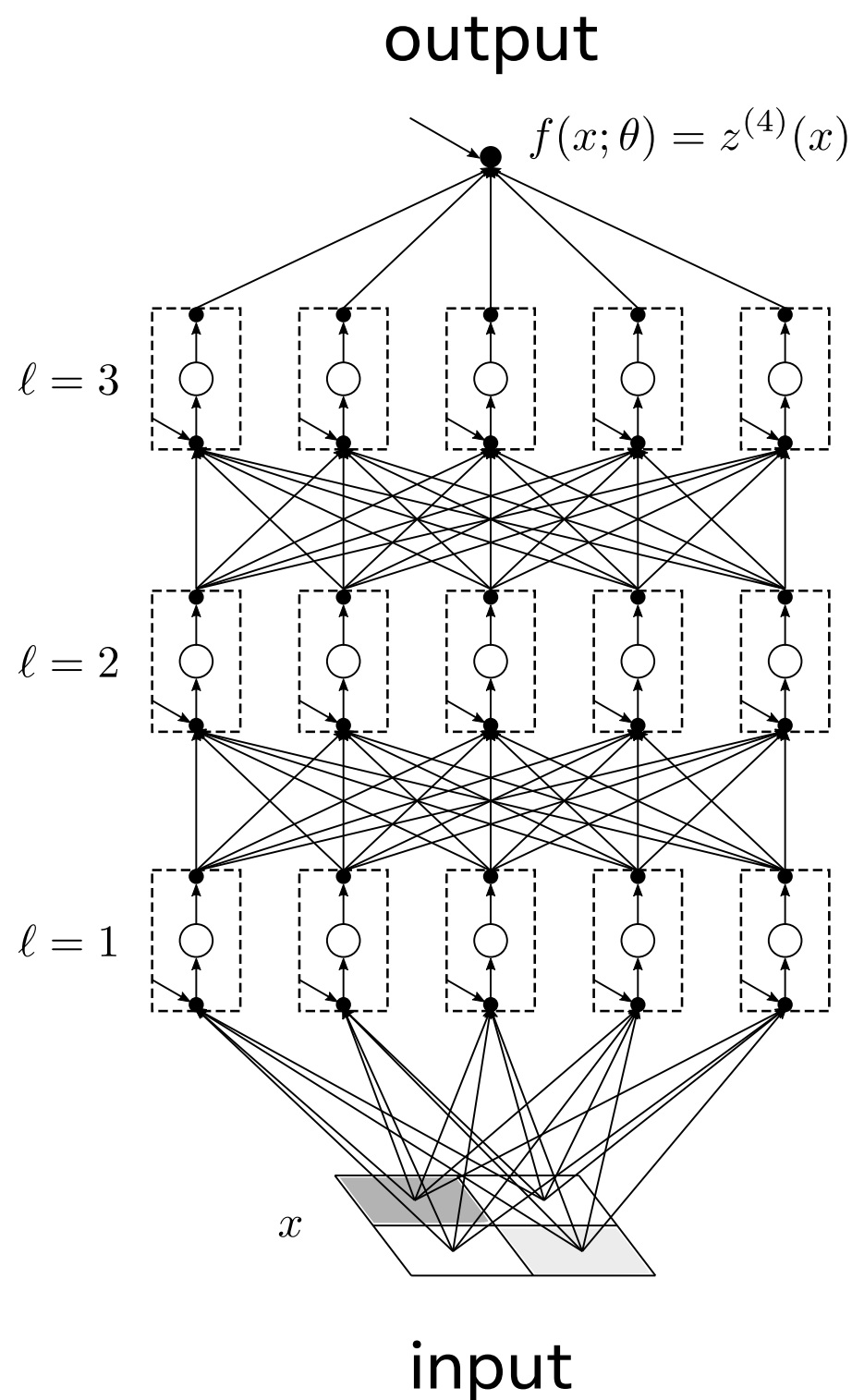
$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)}(x)) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$f(x; \theta) = z^{(L)}(x)$$

- Model parameters:  $\theta_{\mu=1, \dots, P} = \left\{ b_i^{(1)}, W_{ij}^{(1)}, b_i^{(2)}, W_{ij}^{(2)}, \dots, b_i^{(L)}, W_{ij}^{(L)} \right\}$

# Neural Networks



- Function:

$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)}(x)) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

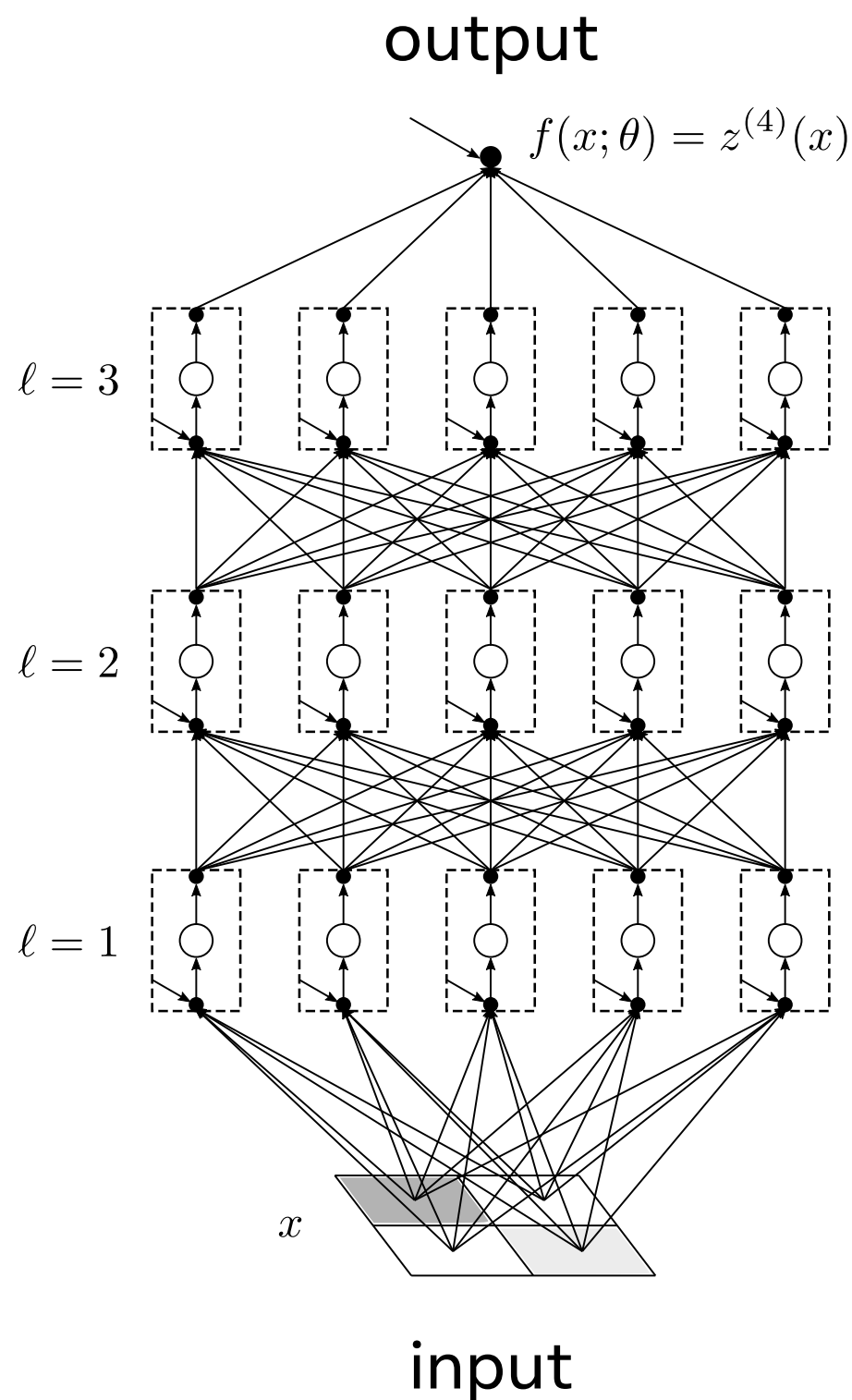
$$f(x; \theta) = z^{(L)}(x)$$

- Model parameters:  $\theta_{\mu=1, \dots, P} = \{b_i^{(1)}, W_{ij}^{(1)}, b_i^{(2)}, W_{ij}^{(2)}, \dots, b_i^{(L)}, W_{ij}^{(L)}\}$

- Initialization distribution: i.i.d. from mean-zero Gaussian with

$$\mathbb{E} \begin{bmatrix} b_{i_1}^{(\ell)} & b_{i_2}^{(\ell)} \end{bmatrix} = \delta_{i_1 i_2} C_b^{(\ell)}, \quad \mathbb{E} \begin{bmatrix} W_{i_1 j_1}^{(\ell)} & W_{i_2 j_2}^{(\ell)} \end{bmatrix} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}$$

# Neural Networks



- Function:

$$z_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$z_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)}(x)) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$f(x; \theta) = z^{(L)}(x)$$

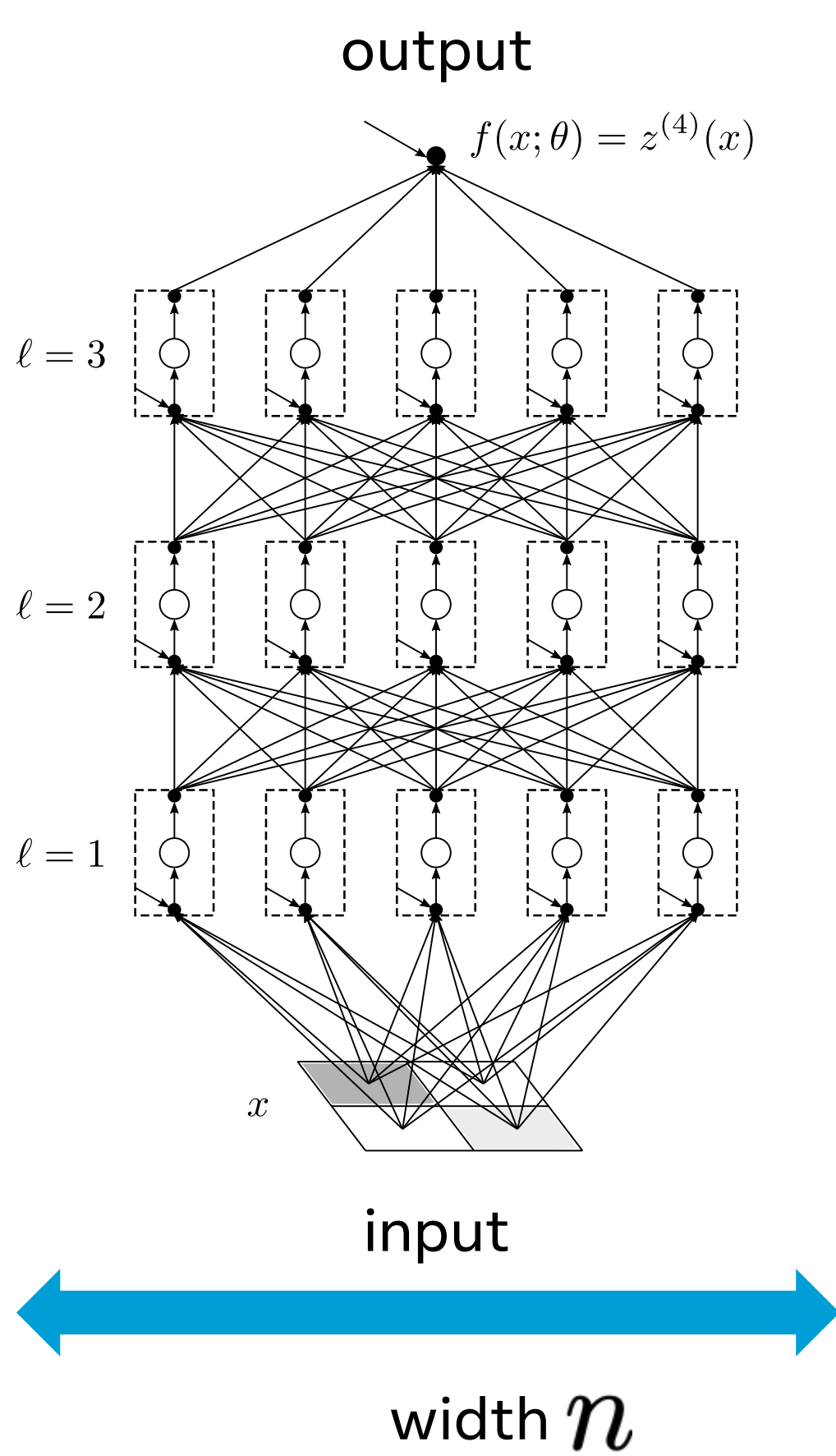
- Model parameters:  $\theta_{\mu=1, \dots, P} = \{b_i^{(1)}, W_{ij}^{(1)}, b_i^{(2)}, W_{ij}^{(2)}, \dots, b_i^{(L)}, W_{ij}^{(L)}\}$

- Initialization distribution: i.i.d. from mean-zero Gaussian with

$$\mathbb{E} \begin{bmatrix} b_{i_1}^{(\ell)} & b_{i_2}^{(\ell)} \end{bmatrix} = \delta_{i_1 i_2} C_b^{(\ell)}, \quad \mathbb{E} \begin{bmatrix} W_{i_1 j_1}^{(\ell)} & W_{i_2 j_2}^{(\ell)} \end{bmatrix} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}$$

good wide limit

# Neural Networks



depth  $L$

# of model parameters

$$P \sim n^2 L$$

# Machine Learning in a Nutshell

- Instantiate a model

$$f_{\text{init}}(x) = f(x; \theta_{\text{init}}) \quad \text{with} \quad \underline{\theta_{\text{init}} \in p(\theta_{\text{init}})}$$

- Train the model, e.g. by gradient descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \right|_{\theta = \theta(t)}$$

- Use the trained model to make predictions

$$p(f_{\text{trained}})$$

mean, variance, etc.



# Problems 1, 2, & 3



# Problems 1, 2, & 3

Trained function, Taylor-expanded around initialization:

$$f_{\text{trained}} = f_{\text{init}} + (\theta_{\text{trained}} - \theta_{\text{init}}) \left. \frac{df}{d\theta} \right|_{\text{init}} + \frac{1}{2} (\theta_{\text{trained}} - \theta_{\text{init}})^2 \left. \frac{d^2 f}{d\theta^2} \right|_{\text{init}} + \dots$$

# Problems 1, 2, & 3

Trained function, Taylor-expanded around initialization:

$$f_{\text{trained}} = f_{\text{init}} + (\theta_{\text{trained}} - \theta_{\text{init}}) \left. \frac{df}{d\theta} \right|_{\text{init}} + \frac{1}{2} (\theta_{\text{trained}} - \theta_{\text{init}})^2 \left. \frac{d^2 f}{d\theta^2} \right|_{\text{init}} + \dots$$

- Problem 1: too many terms in general

# Problems 1, 2, & 3

Trained function, Taylor-expanded around initialization:

$$f_{\text{trained}} = f_{\text{init}} + (\theta_{\text{trained}} - \theta_{\text{init}}) \left. \frac{df}{d\theta} \right|_{\text{init}} + \frac{1}{2} (\theta_{\text{trained}} - \theta_{\text{init}})^2 \left. \frac{d^2 f}{d\theta^2} \right|_{\text{init}} + \dots$$

- Problem 1: too many terms in general
- Problem 2: complicated mapping

$$p(\theta_{\text{init}}) \rightarrow p \left( \theta_{\text{init}}, f_{\text{init}}, \left. \frac{df}{d\theta} \right|_{\text{init}}, \left. \frac{d^2 f}{d\theta^2} \right|_{\text{init}}, \dots \right)$$

statistics at *initialization*

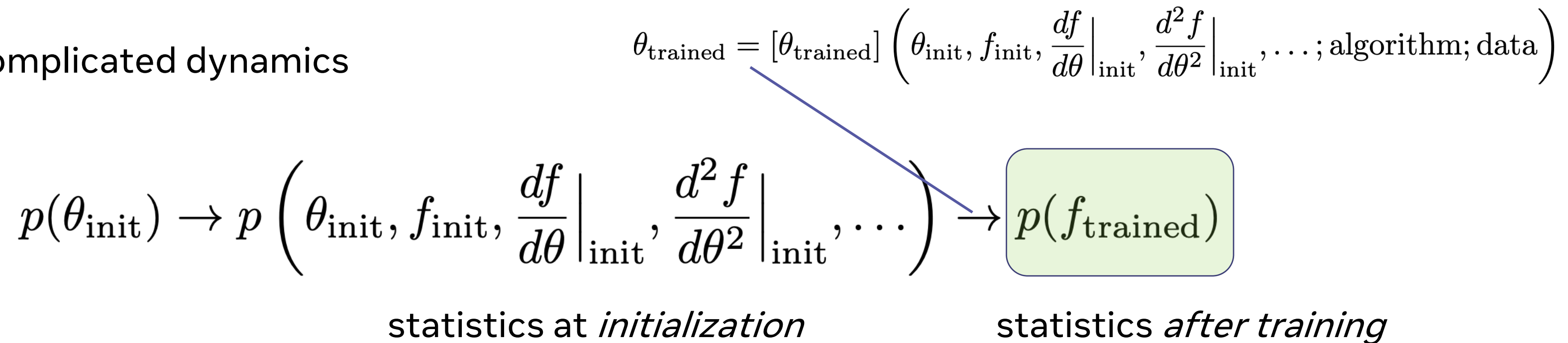


# Problems 1, 2, & 3

Trained function, Taylor-expanded around initialization:

$$f_{\text{trained}} = f_{\text{init}} + (\theta_{\text{trained}} - \theta_{\text{init}}) \left. \frac{df}{d\theta} \right|_{\text{init}} + \frac{1}{2} (\theta_{\text{trained}} - \theta_{\text{init}})^2 \left. \frac{d^2 f}{d\theta^2} \right|_{\text{init}} + \dots$$

- Problem 1: too many terms in general
- Problem 2: complicated mapping
- Problem 3: complicated dynamics



# Despair & Hope

# Despair & Hope

- Microscopic perspective (focusing on individuals):  
the more model parameters, the more complex. We are doomed...

# Despair & Hope

- Microscopic perspective (focusing on individuals):  
the more model parameters, the more complex. We are doomed...
- Macroscopic perspective (focusing on averages):  
the more model parameters, the simpler. We can do this!



# Despair & Hope

- Microscopic perspective (focusing on individuals):  
the more model parameters, the more complex. We are doomed...
- Macroscopic perspective (focusing on averages):  
the more model parameters, the simpler. We can do this!

Simplification when there are infinitely-many neurons in hidden layers  
(a.k.a. law of large numbers)

AND

systematically going beyond that idealized limit  
(a.k.a. perturbation theory)

# Despair & Hope

- Microscopic perspective (focusing on individuals):  
the more model parameters, the more complex. We are doomed...
- Macroscopic perspective (focusing on averages):  
the more model parameters, the simpler. We can do this!

Simplification when there are infinitely-many neurons in hidden layers  
(a.k.a. law of large numbers)

AND

systematically going beyond that idealized limit  
(a.k.a. perturbation theory)

$$p(\theta_{\text{init}}) \rightarrow p\left(\theta_{\text{init}}, f_{\text{init}}, \left.\frac{df}{d\theta}\right|_{\text{init}}, \left.\frac{d^2f}{d\theta^2}\right|_{\text{init}}, \dots\right) \rightarrow p(f_{\text{trained}})$$

---

Statistics become *sparse* & dynamics can be truncated

Trailer for part 2:  $n = \infty$

# Trailer for part 2: $n = \infty$

- Gaussian Prior (& Posterior) [R. Neal (1996), J. Lee+Y. Bahri et al. (ICLR 2018), A. Matthews et al. (ICLR2018)]
- (Neural Tangent) Kernel Learning [A. Jacot, F. Gabriel, & C. Hongler (NeurIPS 2018)]

As simple as we can imagine, but...

# Trailer for part 2: $n = \infty$

- Gaussian Prior (& Posterior) [R. Neal (1996), J. Lee+Y. Bahri et al. (ICLR 2018), A. Matthews et al. (ICLR2018)]
- (Neural Tangent) Kernel Learning [A. Jacot, F. Gabriel, & C. Hongler (NeurIPS 2018)]

As simple as we can imagine, but...

- No Representation Learning(\*)
- No Algorithm Dependence

too simple to describe real deep neural networks

(\*Representation Learning:  
the ability of a model to learn useful representations from data)

# Trailer for part 2: $n = \infty$

- Gaussian Prior (& Posterior) [R. Neal (1996), J. Lee+Y. Bahri et al. (ICLR 2018), A. Matthews et al. (ICLR2018)]
- (Neural Tangent) Kernel Learning [A. Jacot, F. Gabriel, & C. Hongler (NeurIPS 2018)]

As simple as we can imagine, but...

- No Representation Learning
- No Algorithm Dependence

too simple to describe real deep neural networks

A useful starting point but not the end of the story

Trailer for part 3:  $n \gg L$

# Trailer for part 3: $n \gg L$

- *Nearly-Gaussian* Prior (& Posterior)    [§4 (& §6) of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165)]
- *Weakly-Nonlinear* Learning Dynamics    [§11 & §∞]

A bit more complex but tractable, and...



# Trailer for part 3: $n \gg L$

- *Nearly-Gaussian* Prior (& Posterior)    [§4 (& §6) of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165)]
- *Weakly-Nonlinear* Learning Dynamics    [§11 & §∞]

A bit more complex but tractable, and...

- Yes Representation Learning     $\propto \frac{L}{n}$
- Yes Algorithm Dependence

complex enough to capture rich phenomenology of real deep neural networks

# Trailer for part 3: $n \gg L$

- *Nearly-Gaussian* Prior (& Posterior)    [§4 (& §6) of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165)]
- *Weakly-Nonlinear* Learning Dynamics    [§11 & §∞]

A bit more complex but tractable, and...

- Yes Representation Learning     $\propto \frac{L}{n}$
- Yes Algorithm Dependence

complex enough to capture rich phenomenology of real deep neural networks

Qualitatively very different from infinite-width limit

## 2. Neural Networks at Infinite Width

# Some Notations

- Initial outputs:  $\widehat{z}_{i;\delta} = z_i^{(L)}(x_\delta; \theta_{\text{init}})$
- Trained outputs:  $z_{i;\delta}^\star = z_i^{(L)}(x_\delta; \theta_{\text{trained}})$ 
  - introduced sample index  $\delta$
  - dropped  $(L)$
  - hatted initial
  - starred trained

# Training Dynamics

Gradient descent:  $\theta_{\mu}(t + 1) = \theta_{\mu}(t) - \eta \frac{\partial \mathcal{L}}{\partial \theta_{\mu}} \Big|_{\theta = \theta(t)}$

# Training Dynamics

Gradient descent:  $\theta_\mu(t+1) = \theta_\mu(t) - \eta \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\mu} \right)$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W}{n_{\ell-1}}$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W}{n_{\ell-1}}$$

good wide limit



# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}}$$

$O(1/n)$   ~~$(+\dots)$~~

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t)$$

$$- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}}$$

$$\equiv H_{ij;\delta\tilde{\alpha}}(t) \quad \text{Neural Tangent Kernel (NTK)}$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) &= z_{i;\delta}(t) \\ &\quad - \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \\ &\quad \equiv H_{ij;\delta\tilde{\alpha}}(t) \quad \text{Neural Tangent Kernel (NTK)} \end{aligned}$$

Similarly:

$$H_{i_1 i_2; \delta_1 \delta_2}(t+1) = H_{i_1 i_2; \delta_1 \delta_2}(t) + O\left(\frac{1}{n}\right)$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) &= z_{i;\delta}(t) \\ &\quad - \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \\ &\equiv H_{ij;\delta\tilde{\alpha}}(t) = \hat{H}_{ij;\delta\tilde{\alpha}} \quad \text{“frozen” NTK} \end{aligned}$$

Similarly:

$$H_{i_1 i_2; \delta_1 \delta_2}(t+1) = H_{i_1 i_2; \delta_1 \delta_2}(t) + \cancel{O\left(\frac{1}{n}\right)}$$

# Training Dynamics

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \hat{H}_{ij;\delta\tilde{\alpha}} \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}}$$

## Solving “Problem 3” (Dynamics)

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \hat{H}_{ij;\delta\tilde{\alpha}} \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}}$$

# Solving “Problem 3” (Dynamics)

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \hat{H}_{ij;\delta\tilde{\alpha}} \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}}$$

“E.g.,” (\*) for  $\mathcal{L} = \frac{1}{2} \sum_{i,\tilde{\alpha}} (z_{i;\tilde{\alpha}} - y_{i;\tilde{\alpha}})^2$

(\*It turns out that the detailed forms of the loss/scheduling won't matter:  
*algorithm independence*, §10.2.2 of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165))

# Solving “Problem 3” (Dynamics)

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \hat{H}_{ij;\delta\tilde{\alpha}} [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}]$$

“E.g.,” (\*) for  $\mathcal{L} = \frac{1}{2} \sum_{i,\tilde{\alpha}} (z_{i;\tilde{\alpha}} - y_{i;\tilde{\alpha}})^2$

(\*It turns out that the detailed forms of the loss/scheduling won't matter:  
*algorithm independence*, §10.2.2 of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165))



# Solving “Problem 3” (Dynamics)

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \hat{H}_{ij;\delta\tilde{\alpha}} [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}]$$

“E.g.,” (\*) for  $\mathcal{L} = \frac{1}{2} \sum_{i,\tilde{\alpha}} (z_{i;\tilde{\alpha}} - y_{i;\tilde{\alpha}})^2$

→ (exponentially)

$$z_{i;\delta}^* = \hat{z}_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\delta\tilde{\alpha}_1} \left( \hat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\hat{z}_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2})$$

(\*It turns out that the detailed forms of the loss/scheduling won't matter:  
*algorithm independence*, §10.2.2 of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165))

# Solving “Problem 3” (Dynamics)

$$z_{i;\delta}^* = \hat{z}_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\delta\tilde{\alpha}_1} \left( \hat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\hat{z}_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2})$$

$$p(\theta_{\text{init}}) \rightarrow p(\hat{z}, \hat{H}) \rightarrow p(z^*)$$

# Solutions to “Problems 1 & 2”

$$p(\theta_{\text{init}}) \rightarrow p(\hat{z}, \hat{H}) \rightarrow p(z^*)$$

# Solutions to “Problems 1 & 2”

- Gaussian distribution [R. Neal (1996), J. Lee+Y. Bahri et al. (ICLR 2018), A. Matthews et al. (ICLR2018)]

$$p(\hat{z}_{i;\delta}) \propto \exp \left[ -\frac{1}{2} \sum_{i,\delta_1,\delta_2} (K^{-1})^{\delta_1\delta_2} \hat{z}_{i;\delta_1} \hat{z}_{i;\delta_2} \right]$$

- Deterministic NTK [A. Jacot, F. Gabriel, & C. Hongler (NeurIPS 2018)]

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2} = \delta_{i_1 i_2} \Theta_{\delta_1 \delta_2}$$

$$p(\theta_{\text{init}}) \rightarrow p(\hat{z}, \hat{H}) \rightarrow p(z^*)$$

# Solved EVERYTHING

$p(z_{i;\delta}^*)$  Gaussian distribution

with mean  $m_{i;\delta} = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \text{train}} \Theta_{\delta \tilde{\alpha}_1} (\tilde{\Theta}^{-1})^{\tilde{\alpha}_1 \tilde{\alpha}_2} y_{i;\tilde{\alpha}_2}$

and variance involving both  $K_{\delta_1 \delta_2}$ ,  $\Theta_{\delta_1 \delta_2}$

# Solved EVERYTHING but...

$p(z_{i;\delta}^*)$  Gaussian distribution

with mean  $m_{i;\delta} = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \text{train}} \Theta_{\delta \tilde{\alpha}_1} (\tilde{\Theta}^{-1})^{\tilde{\alpha}_1 \tilde{\alpha}_2} y_{i;\tilde{\alpha}_2}$

and variance involving both  $K_{\delta_1 \delta_2}, \Theta_{\delta_1 \delta_2}$

- No Representation Learning (linear model with random features)
- No Algorithm Dependence (GD, Newton, SGD with decreasing learning rate, ...)

too simple to describe real deep neural networks

# 3. Neural Networks at Finite Width

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) \quad \text{NTK } H(t) \\ - \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \\ + \dots$$



# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) \quad \text{NTK } H(t)$$

$$- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK) } dH(t)$$

$$+ \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left( \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1; \tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2; \tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1; \tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2; \tilde{\alpha}_2}}$$

$$+ \dots$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned}
 z_{i;\delta}(t+1) &= z_{i;\delta}(t) && \text{NTK } H(t) \\
 &- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} && \text{differential of NTK (dNTK) } dH(t) \\
 &+ \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \left( \sum_{\mu_1,\nu_1,\mu_2,\nu_2} \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \\
 &- \frac{\eta^3}{6} \sum \left( \sum_{\mu_1,\nu_1,\mu_2,\nu_2,\mu_3,\nu_3} \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \lambda_{\mu_3\nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}} \\
 &+ \dots && \text{ddNTK } ddH(t)
 \end{aligned}$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) \quad \text{NTK } H(t)$$

$$- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK) } dH(t)$$

$$O(1/n) \left( + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left( \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \right.$$

$$- \frac{\eta^3}{6} \sum \left( \sum \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \lambda_{\mu_3 \nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}} \quad \text{ddNTK } ddH(t)$$

$$O(1/n^2) \left( + \dots \right)$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) \quad \text{NTK } H(t)$$

$$- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK) } dH(t)$$

$$O(1/n) \left( + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left( \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \right.$$

$$- \frac{\eta^3}{6} \sum \left( \sum \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \lambda_{\mu_3 \nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}}$$

$$\left. + \dots \right) \quad \text{ddNTK } ddH(t)$$

# Training Dynamics

Gradient descent: 
$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_{\nu=1}^P \lambda_{\mu\nu} \left( \sum_{\tilde{\alpha} \in \mathcal{B}_{\text{train}}} \sum_j \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) \quad \text{NTK } H(t)$$

$$- \eta \sum_{j,\tilde{\alpha}} \left( \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK) } dH(t)$$

$$O(1/n) \left( + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left( \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \right.$$

$$\left. - \frac{\eta^3}{6} \sum \left( \sum \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \lambda_{\mu_3 \nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}} \right.$$

$$\left. + \dots \right) \quad \text{ddNTK } ddH(t)$$

Similarly some dynamical equations for NTK and dNTK (while ddNTK is frozen at this order)

# Solving “Problem 3” (Dynamics)

[...long song & dance with dynamical perturbation theory to get  $z^* \left( \hat{z}, \hat{H}, \widehat{dH}, \widehat{ddH} \right) \dots]$

## Solving “Problem 3” (Dynamics)

$$z_{i;\delta}^* = \hat{z}_{i;\delta} - \sum \hat{H}_{ij;\delta\tilde{\alpha}_1} \left( \hat{H}^{-1} \right)^{jk;\tilde{\alpha}_1\tilde{\alpha}_2} [\hat{z}_{k;\tilde{\alpha}} - y_{k;\tilde{\alpha}}] \\ + \text{despicable}(y, \hat{z}, \hat{H}, \widehat{dH}, \widehat{ddH}; \text{algorithm})$$

$$H^* \neq \hat{H}$$

## Solving “Problem 3” (Dynamics)

$$z_{i;\delta}^* = \hat{z}_{i;\delta} - \sum \hat{H}_{ij;\delta\tilde{\alpha}_1} \left( \hat{H}^{-1} \right)^{jk;\tilde{\alpha}_1\tilde{\alpha}_2} [\hat{z}_{k;\tilde{\alpha}} - y_{k;\tilde{\alpha}}] \\ + \text{despicable}(y, \hat{z}, \hat{H}, \widehat{dH}, \widehat{ddH}; \text{algorithm})$$

$$H^* \neq \hat{H}$$

- NTK now “defrosted”: indicative of Representation Learning!
- Solution depends on the algorithm: Algorithm Dependence!



# Solving “Problem 3” (Dynamics)

$$\begin{aligned}
 & z_{i;\delta}(t = \infty) \tag{\infty.141} \\
 & \equiv z_{i;\delta}^F(t = \infty) + z_{i;\delta}^I(t = \infty) \\
 & = z_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{H}_{ij;\delta\tilde{\alpha}_1} \left( \widehat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & \quad + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] \\
 & \quad \quad \quad \times Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & \quad + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] \\
 & \quad \quad \quad \times Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\delta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_8\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + O\left(\frac{1}{n^2}\right).
 \end{aligned}$$

# Solving “Problem 3” (Dynamics)

$$\begin{aligned}
 & z_{i;\delta}(t = \infty) && (\infty.141) \\
 & \equiv z_{i;\delta}^F(t = \infty) + z_{i;\delta}^I(t = \infty) \\
 & = z_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{H}_{ij;\delta\tilde{\alpha}_1} \left( \widehat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] \\
 & \quad \times Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] \\
 & \quad \times Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \times Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \times Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\delta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_8\tilde{\alpha}_3} \right] \\
 & \quad \times Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \times Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + O\left(\frac{1}{n^2}\right).
 \end{aligned}$$

algorithm projectors

# Solving “Problem 3” (Dynamics)

$$\begin{aligned}
 & z_{i;\delta}(t = \infty) \tag{\infty.141} \\
 & \equiv z_{i;\delta}^F(t = \infty) + z_{i;\delta}^I(t = \infty) \\
 & = z_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{H}_{ij;\delta\tilde{\alpha}_1} \left( \widehat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & \quad + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] \\
 & \quad \quad \quad \times Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & \quad + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[ \widehat{dH}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\delta\tilde{\alpha}_5} \widetilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] \\
 & \quad \quad \quad \times Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_I H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\delta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_8\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & \quad + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[ \widehat{dd_{II} H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\delta\tilde{\alpha}_7} \widetilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd_{II} H}_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] \\
 & \quad \quad \quad \times Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6})
 \end{aligned}$$

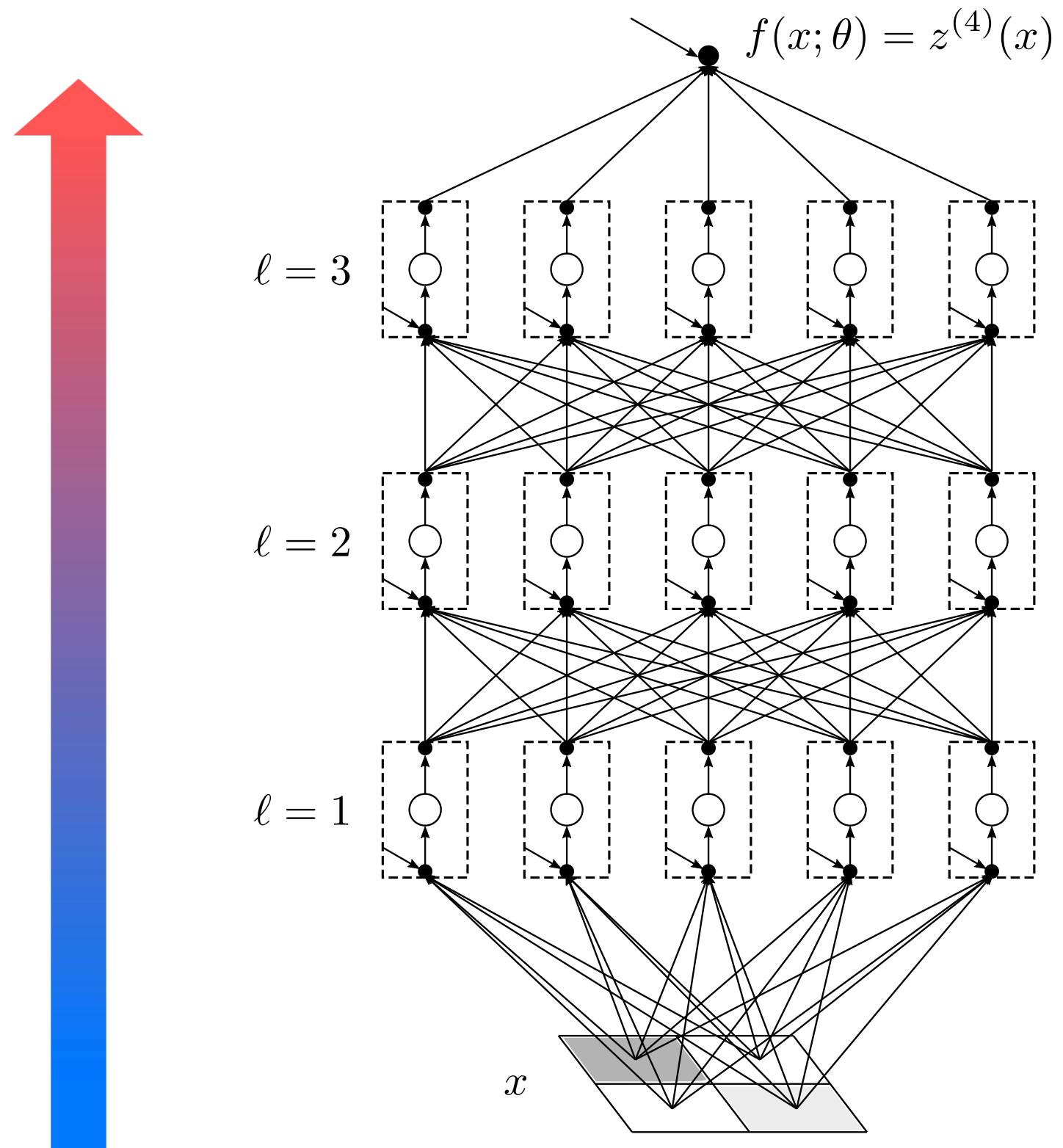
$$p(\theta_{\text{init}}) \rightarrow p(\hat{z}, \widehat{H}, \widehat{dH}, \widehat{ddH}) \xrightarrow{z^*} p(z^*)$$

# Solutions to “Problems 1 & 2”

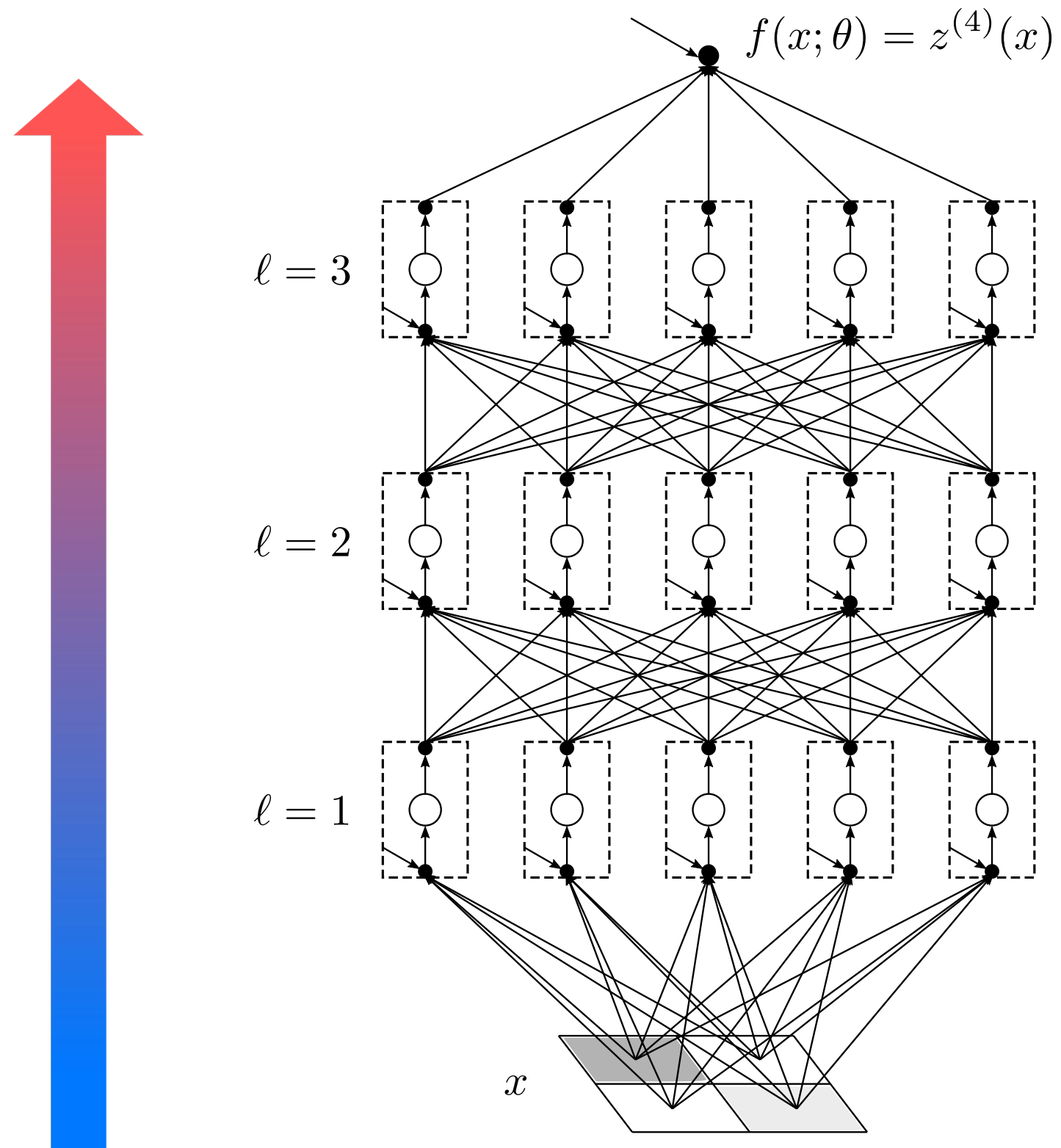
[See: §4, §8, §11.2, & § $\infty$ .3 of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165)]

$$p(\theta_{\text{init}}) \xrightarrow{\circlearrowright} p(\hat{z}, \hat{H}, \widehat{dH}, \widehat{ddH}) \rightarrow p(z^*)$$

# Strategy: Always Forward

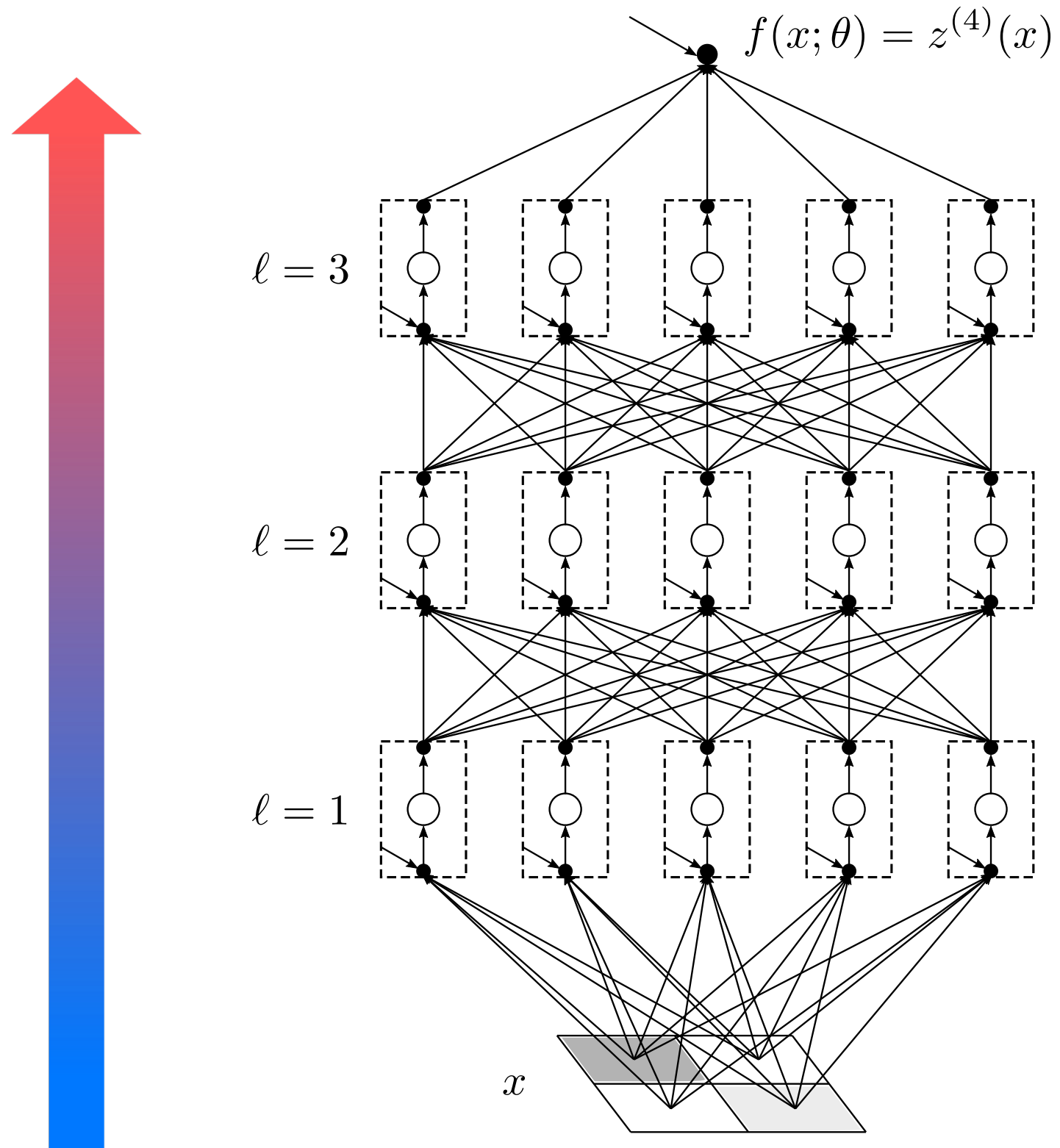


# Strategy: Always Forward



$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$

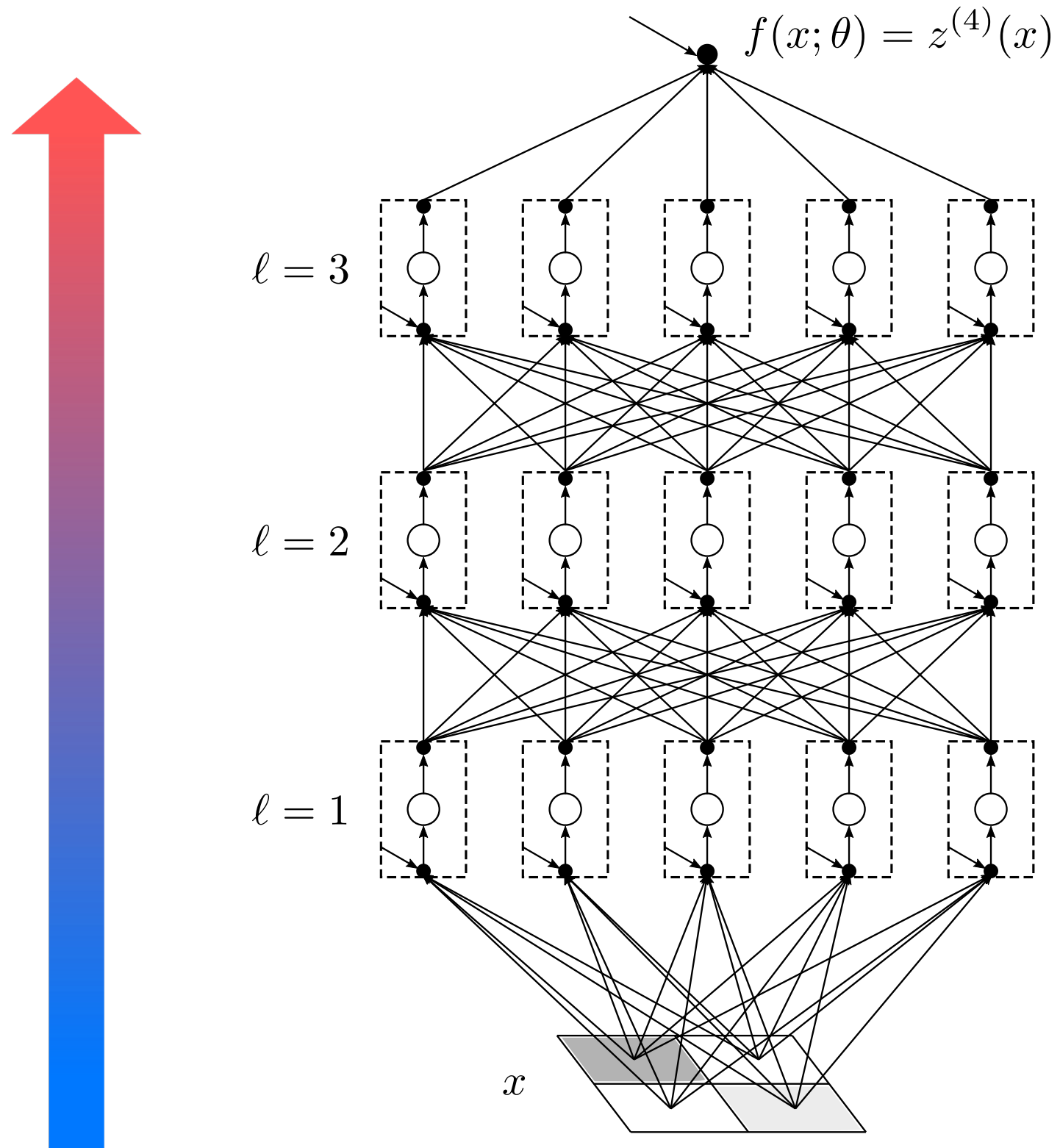
# Strategy: Always Forward



$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{dH}^{(2)}, \dots)$$

$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$

# Strategy: Always Forward



$$p(\hat{z}^{(4)}, \hat{H}^{(4)}, \widehat{dH}^{(4)}, \dots)$$

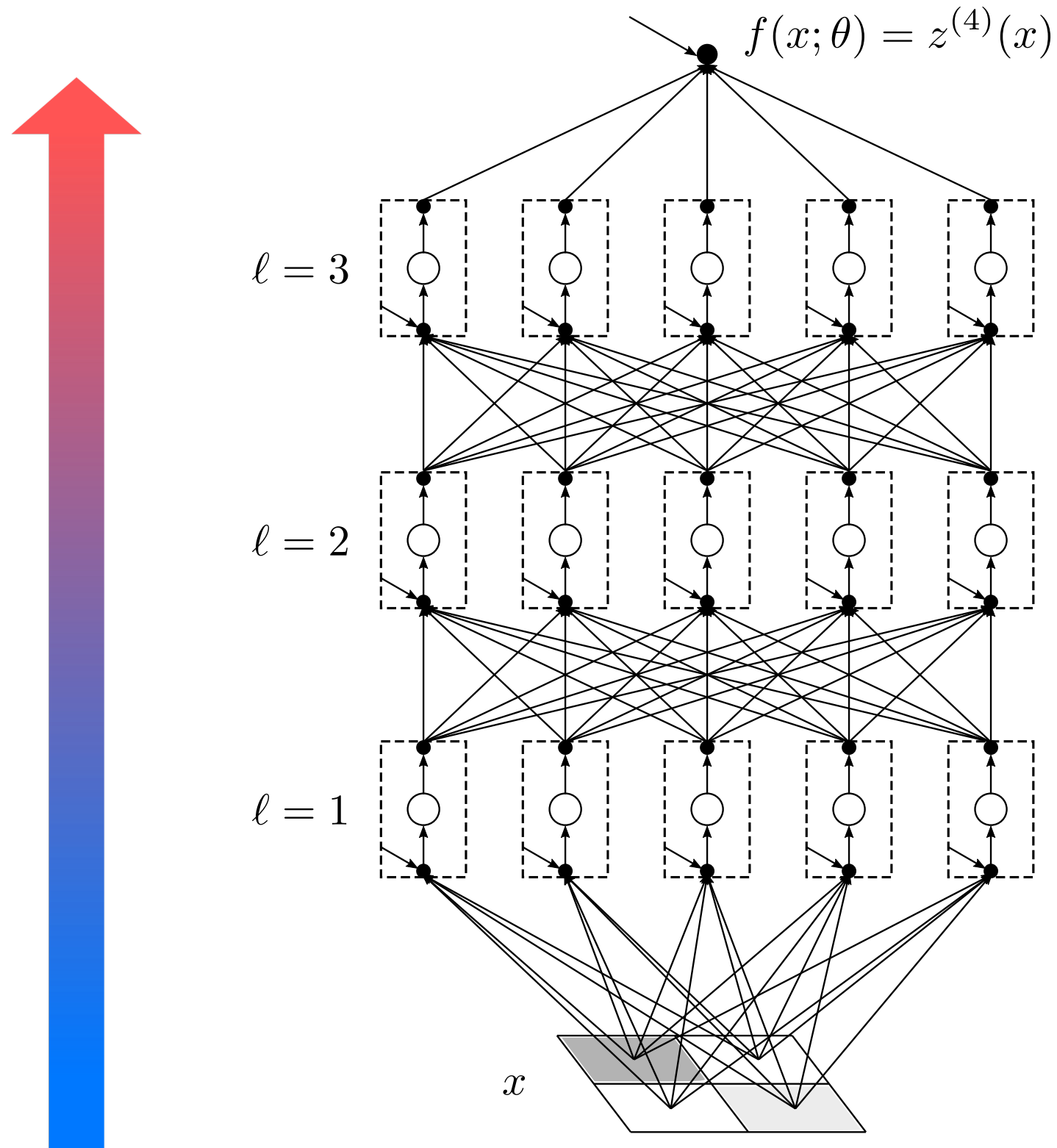
$$p(\hat{z}^{(3)}, \hat{H}^{(3)}, \widehat{dH}^{(3)}, \dots)$$

$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{dH}^{(2)}, \dots)$$

$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$



# Strategy: Always Forward



$$p(\hat{z}^{(4)}, \hat{H}^{(4)}, \widehat{dH}^{(4)}, \dots)$$

$$p(\hat{z}^{(3)}, \hat{H}^{(3)}, \widehat{dH}^{(3)}, \dots)$$

$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{dH}^{(2)}, \dots)$$

$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$

(RG-flow interpretation: §4.6 of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165))

# Solved EVERYTHING

$p(z_i^*; \delta)$  nearly-Gaussian

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$

# Solved EVERYTHING

$p(z_i^*; \delta)$  nearly-Gaussian

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$

non-Gaussianity

# Solved EVERYTHING

$p(z_i^*; \delta)$  nearly-Gaussian

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$

non-Gaussianity

NTK mean

# Solved EVERYTHING

$p(z_i^*; \delta)$  nearly-Gaussian

NTK fluctuations (agitated NTK)

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$

non-Gaussianity

NTK mean

# Solved EVERYTHING

$$p(z_i^*; \delta)$$

nearly-Gaussian

NTK fluctuations (agitated NTK)

ddNTK (defrosted NTK)

$$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$$

non-Gaussianity

NTK mean

dNTK (defrosted NTK)

# Solved EVERYTHING

$$p(z_i^*; \delta)$$

nearly-Gaussian

NTK fluctuations (agitated NTK)

ddNTK (defrosted NTK)

$$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$$

non-Gaussianity

NTK mean

dNTK (defrosted NTK)

\* all recursively computable

# Solved EVERYTHING and...

$p(z_i^*; \delta)$  nearly-Gaussian

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, \underbrace{P_{\dots}, Q_{\dots}}_{\text{dNTK}}, \overbrace{R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}}^{\text{ddNTK}}$

- Yes Representation Learning (cubic model with evolving features)
- Yes Algorithm Dependence (encapsulated by algorithm projectors)

$$\propto \frac{L}{n}$$

complex enough to capture rich phenomenology of real deep neural networks



# A Word about Overly-Deep Neural Networks

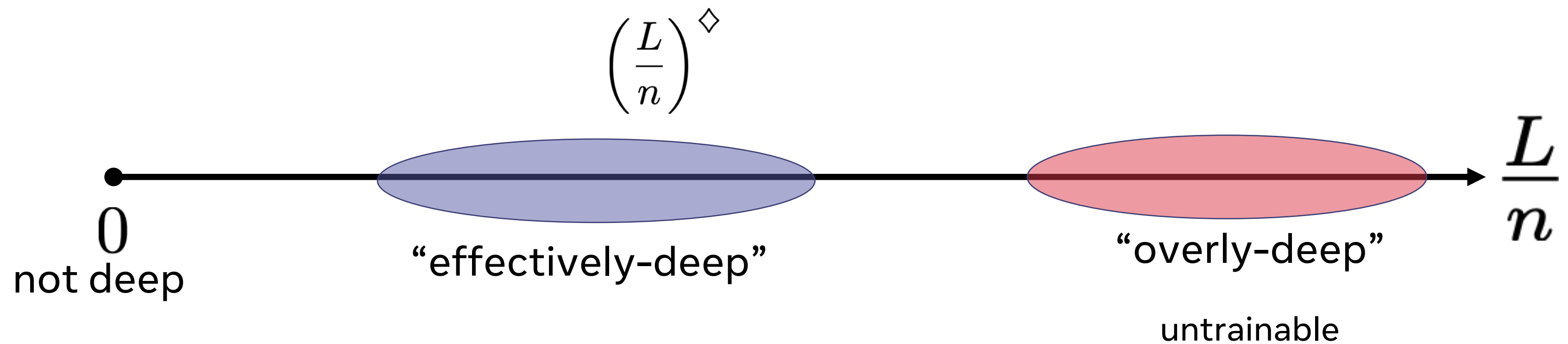
$p(z_i^*; \delta)$  nearly-Gaussian

NTK fluctuations

$G_{\delta_1 \delta_2}, V_{(\delta_1 \delta_2)(\delta_3 \delta_4)}, H_{\delta_1 \delta_2}, A_{\delta_1 \delta_2 \delta_3 \delta_4}, B_{\dots}, D_{\dots}, F_{\dots}, P_{\dots}, Q_{\dots}, R_{\dots}, S_{\dots}, T_{\dots}, U_{\dots}$

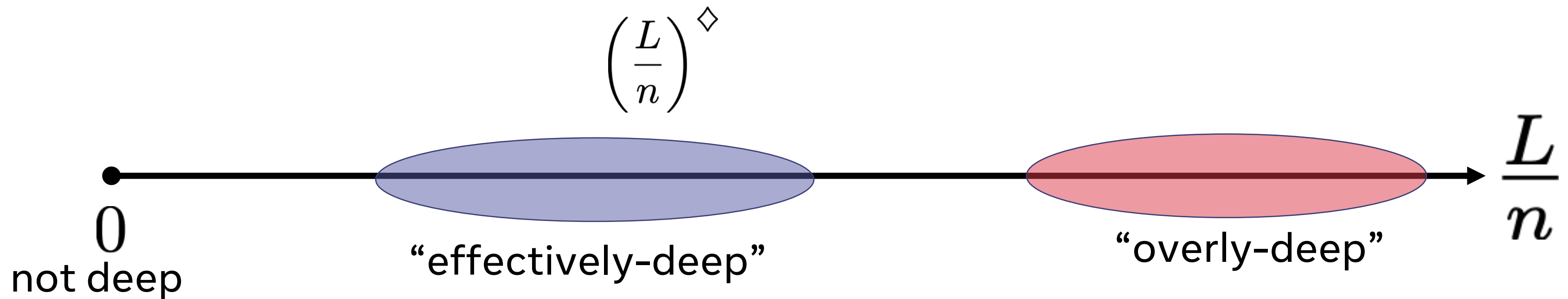
- Instantiation-to-instantiation fluctuations  $\propto \frac{L}{n}$

# A Word about Overly-Deep Neural Networks



# Summary

- $n = \infty$   
simple, no representation learning, no algorithm dependence
- $n \gg L$   
a little more complex but tractable, yes representation learning & algorithm dependence
- $L \gg n$   
too complex, chaotic, untrainable



# 4. The Principles

# The Principle of Sparsity for WIDE Neural Networks

$$p(\theta) \xrightarrow{\text{blue}} p\left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right) \xrightarrow{\text{green}} p(z^*)$$

$z^* \left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right)$

statistics at *initialization*                      statistics *after training*

# The Principle of Sparsity for WIDE Neural Networks

$$p(\theta) \xrightarrow{\text{blue}} p\left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right) \xrightarrow{\text{green}} p(z^*)$$

$z^* \left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right)$

statistics at *initialization*                      statistics *after training*

- Infinite width:

$$p\left(\hat{z}, \hat{H}\right) \text{ specified by } G^{(L)}, H^{(L)} ; \text{linear dynamics}$$

# The Principle of Sparsity for WIDE Neural Networks

$$p(\theta) \xrightarrow{\text{blue circle}} p\left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right) \xrightarrow{\text{green circle}} p(z^*)$$

$z^* \left(\hat{z}, \hat{H}, \widehat{dH}, \dots\right)$   
 statistics at *initialization*                      statistics *after training*

- Infinite width:

$$p\left(\hat{z}, \hat{H}\right) \text{ specified by } G^{(L)}, H^{(L)} ; \text{linear dynamics}$$

- Large-but-finite width at  $O\left(\frac{L}{n}\right) [n_1, n_2, \dots, n_{L-1} \gg L]$ :

$$p\left(\hat{z}, \hat{H}, \widehat{dH}, \widehat{ddH}\right) \text{ specified by}$$

$$G^{(L)}, H^{(L)}, V^{(L)}, A^{(L)}, B^{(L)}, D^{(L)}, F^{(L)}, P^{(L)}, Q^{(L)}, R^{(L)}, S^{(L)}, T^{(L)}, U^{(L)} ; \text{cubic dynamics}$$

# The Principle of Criticality for DEEP Neural Networks

$(C_b, C_W)$  critical

optimal initialization hyperparameters for deep neural networks



# The Principle of Criticality for DEEP Neural Networks

- Taming *exponential* exploding/vanishing kernel problem:  
Poole et al. (NeurIPS2016); Raghu et al. (ICML2016); Schoenholz et al. (ICLR2017); §3 (DLN)+§5 (general) of [arXiv:2106.10165](https://arxiv.org/abs/2106.10165)
- Taming *exponential* exploding/vanishing gradient problem: §9.4
- Bayesian evidence: §6.3.1
- Generalization error: §10.3
- Mutual information: §A.2



$(C_b, C_W)$  critical

optimal initialization hyperparameters for deep neural networks

# The Principle of (Layer) Equivalence

- Taming *polynomial* exploding/vanishing gradient problem: §9.4



How to scale learning rates with depth  
such that all groups of model parameters contribute equally

# The Principle of Typicality

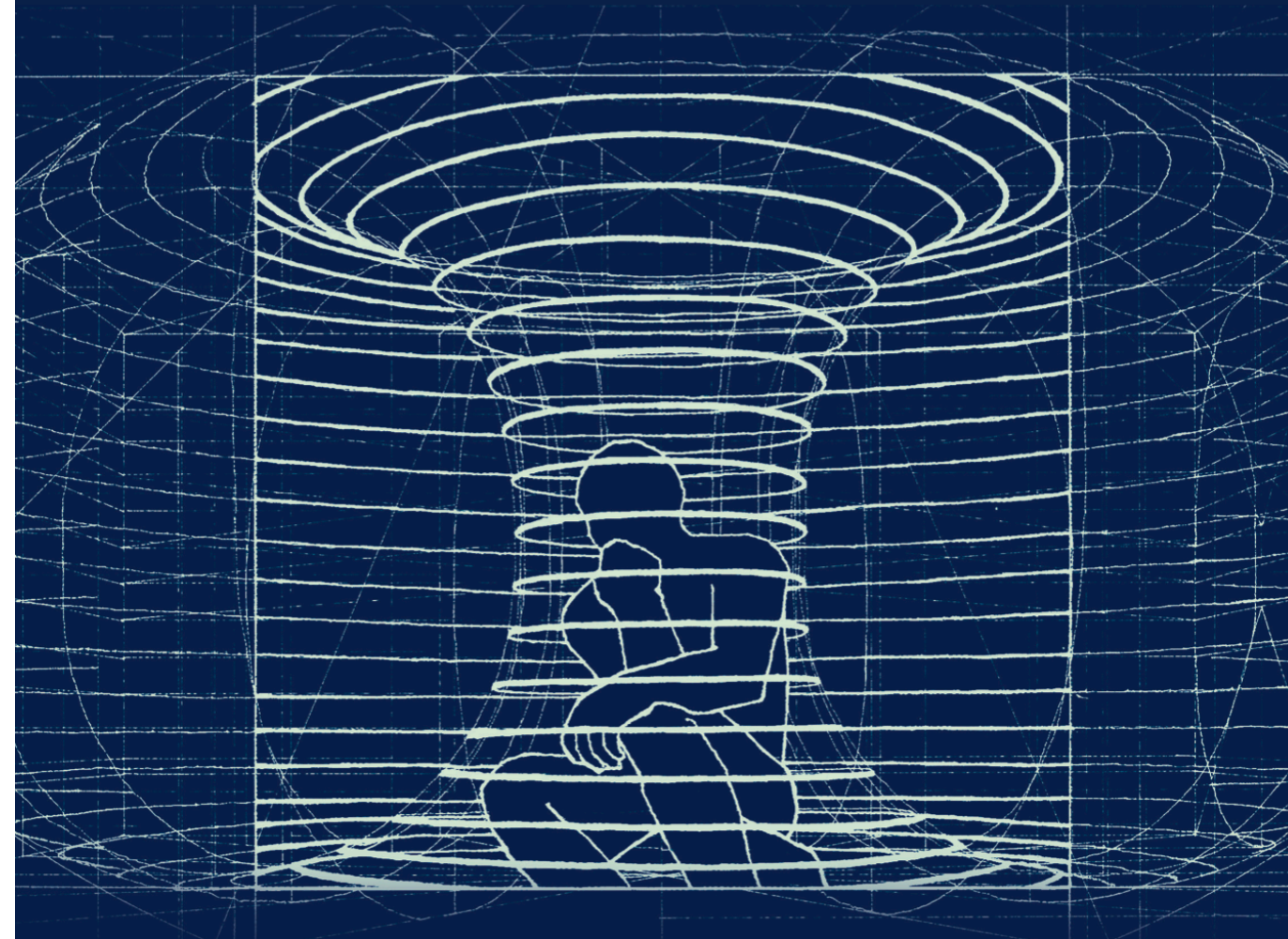
The way we study things:  
analyze statistics and ask what happens *typically*

# The Principle of Universality

The study of criticality can be organized into  
various universality classes of activation functions

# THE PRINCIPLES OF DEEP LEARNING THEORY

An Effective Theory Approach  
to Understanding Neural Networks



Daniel A. Roberts and Sho Yaida  
based on research in collaboration with Boris Hanin

[arXiv:2106.10165](https://arxiv.org/abs/2106.10165)