neural networks

field theories

# Machine Learning and (Large-N) Field Theory

Zhengkang "Kevin" Zhang (University of Utah)

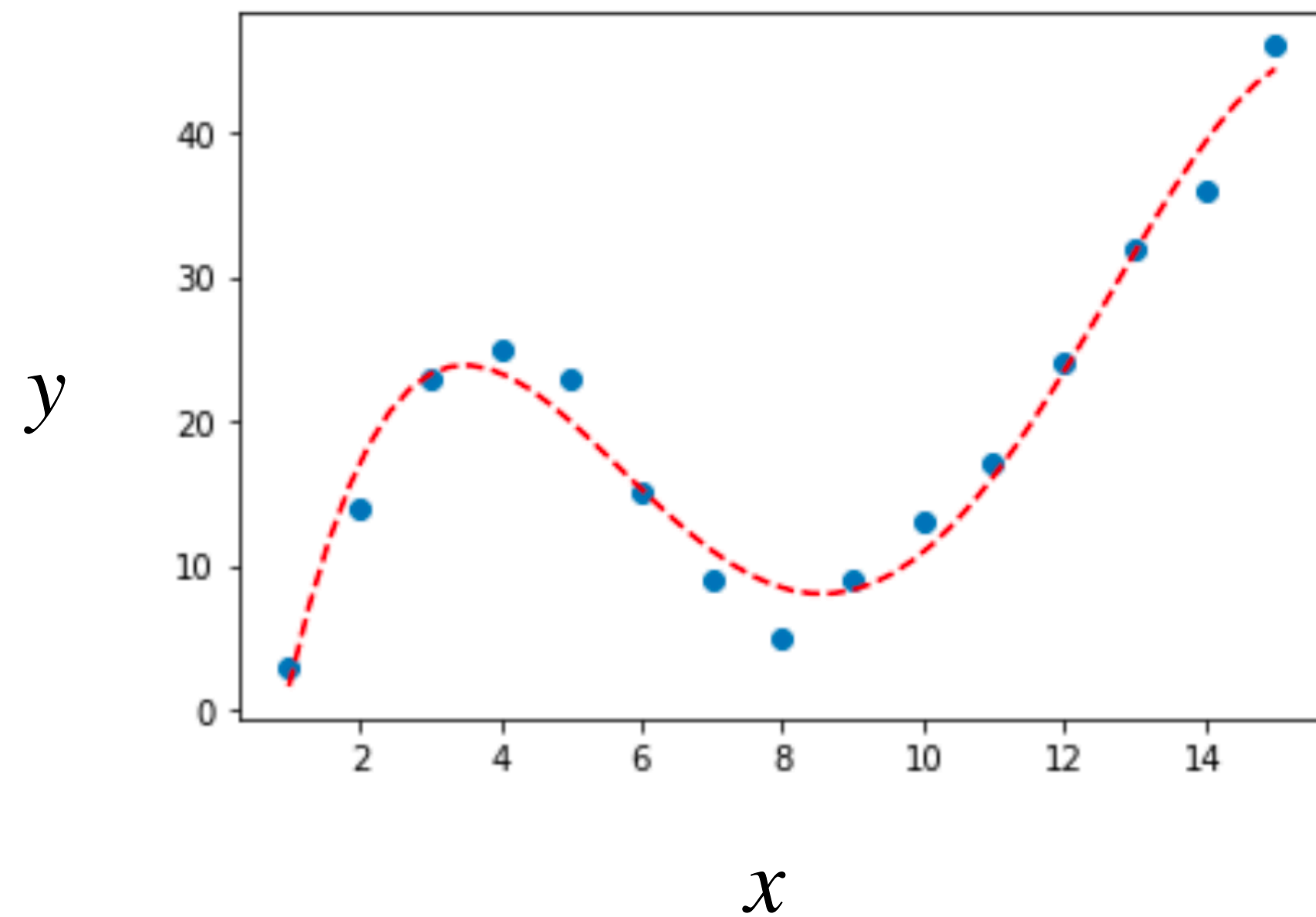# What is a (deep) neural network?

Goal (supervised learning): learn a function $y = f(\vec{x})$ from training dataset $(\vec{x}_\alpha, y_\alpha)$.



| $x$ Image | $y$ Label |
|-----------|-----------|
| | Cat |
| | Cat |
| | Dog |
| | Dog |

# What is a (deep) neural network?

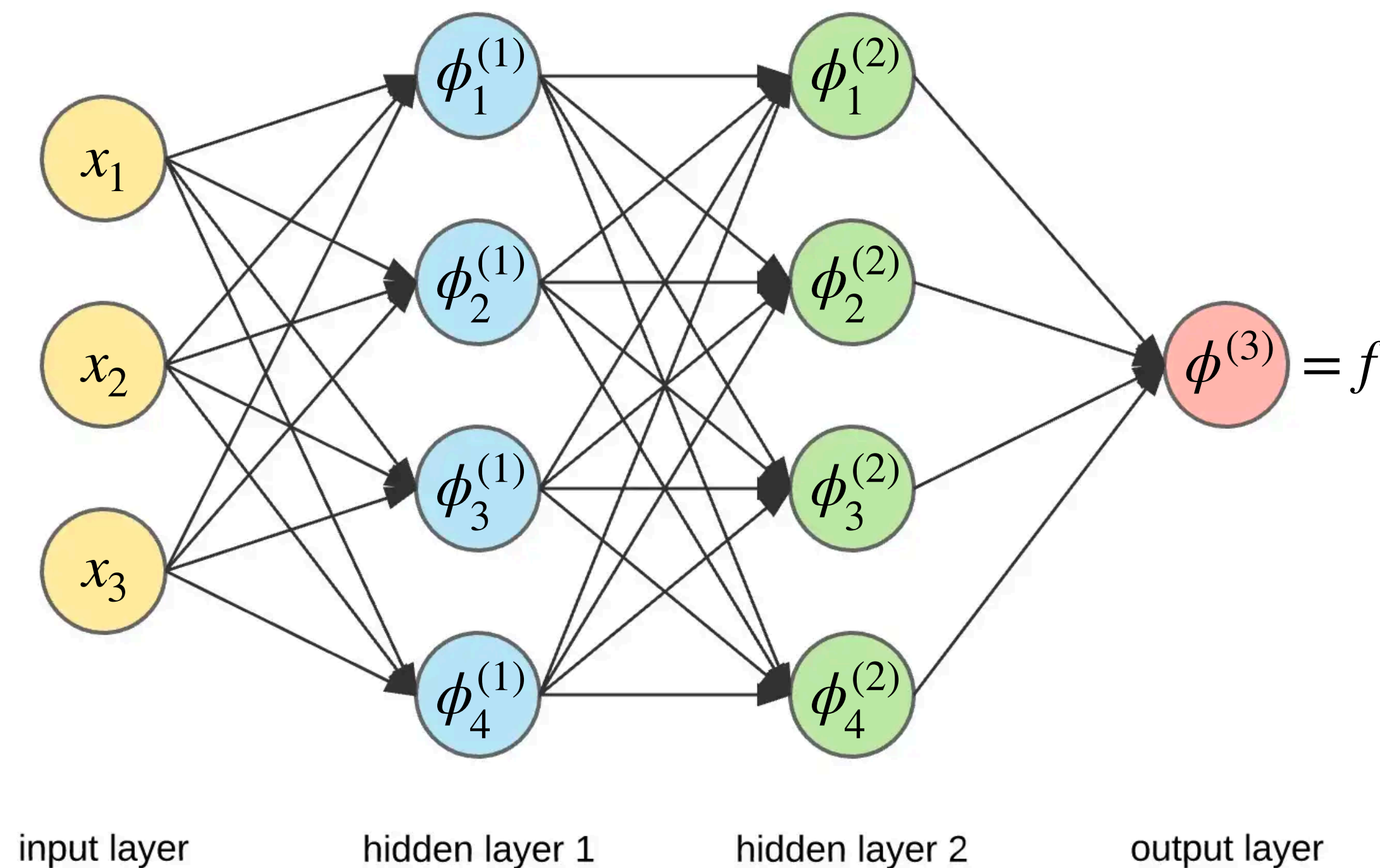Goal (supervised learning): learn a function $y = f(\vec{x})$ from training dataset $(\vec{x}_\alpha, y_\alpha)$.

Neural network = parameterized function with a huge number of parameters ("expressive" enough to represent complicated functions).



input layer     hidden layer 1     hidden layer 2     output layer

$$\phi_i^{(1)}(\vec{x}) = \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j + b_i^{(1)},$$

nonlinear activation function (e.g. tanh)

$$\phi_i^{(\ell)}(\vec{x}) = \sum_{j=1}^{n_{\ell-1}} W_{ij}^{(\ell)} \sigma\left(\phi_j^{(\ell-1)}(\vec{x})\right) + b_i^{(\ell)} \quad (\ell \geq 2).$$

weights      biases

trainable parameters:

- randomly initialized
- then updated to fit training data

# Neural networks ↔ field theories

Ensemble of networks, randomly initialized.

Neurons ↔ scalar fields $\phi(\vec{x})$ .

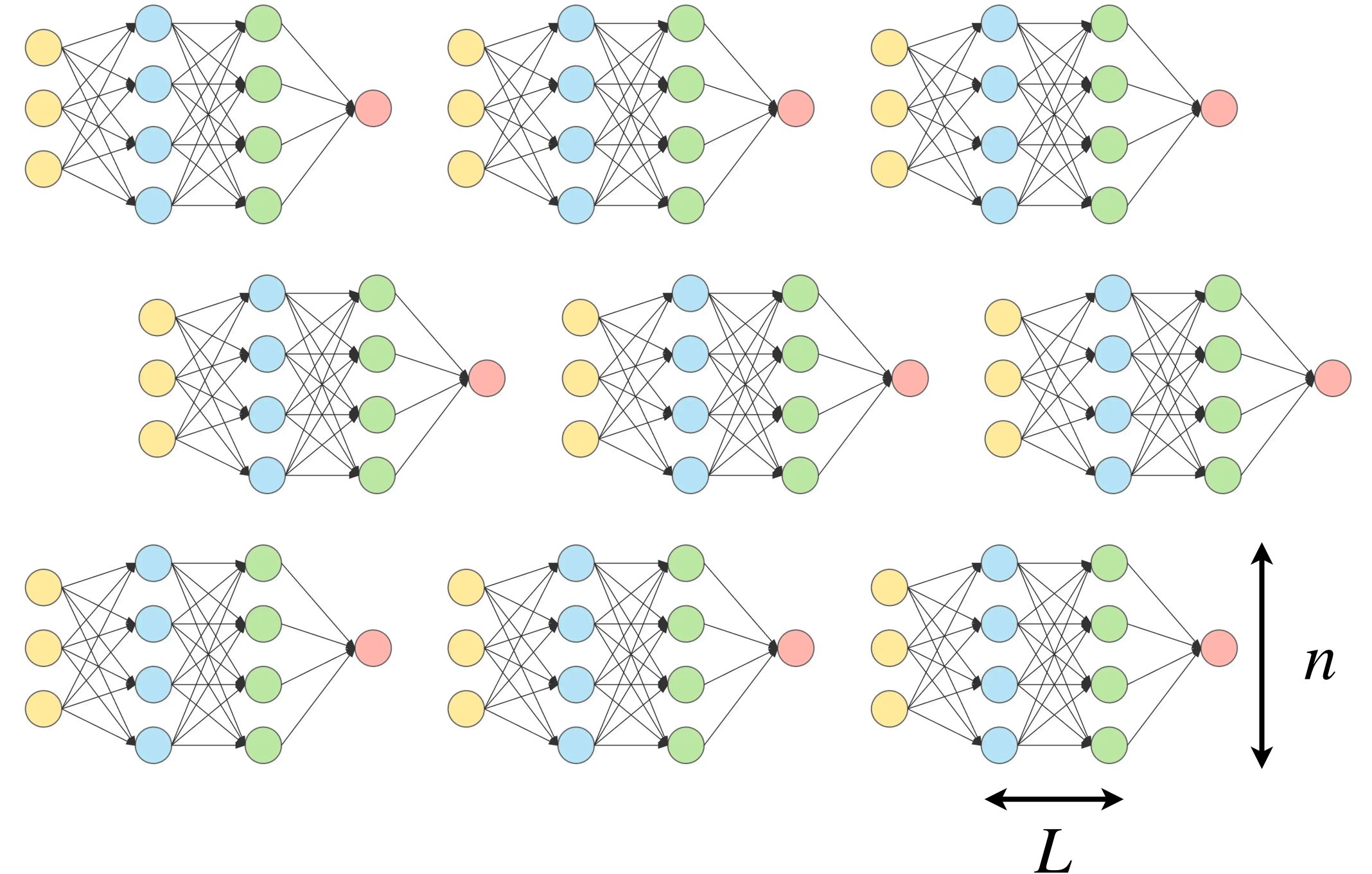Ensemble statistics ↔ action: $P(\phi) = e^{-S[\phi]}$ .

$$\left\langle \phi_{i_1}^{(\ell)}(\vec{x}_1) \, \ldots \, \phi_{i_{2k}}^{(\ell)}(\vec{x}_{2k}) \right\rangle = \int \mathcal{D}\phi \, \phi_{i_1}^{(\ell)}(\vec{x}_1) \, \ldots \, \phi_{i_{2k}}^{(\ell)}(\vec{x}_{2k}) \, e^{-S[\phi]}$$

Evolution with layer $\ell$ ↔ RG flow.

Infinitely-wide networks* $(n \to \infty)$ ↔ free theories.      * Neal '96. Williams '96.

Wide networks $(n \gg L)$ ↔ weakly-interacting theories (perturbative).

4

# Diagrammatic framework

In addition to reproducing known results for lower-point correlators, we were able to push 1/n calculations to higher orders using diagrams.
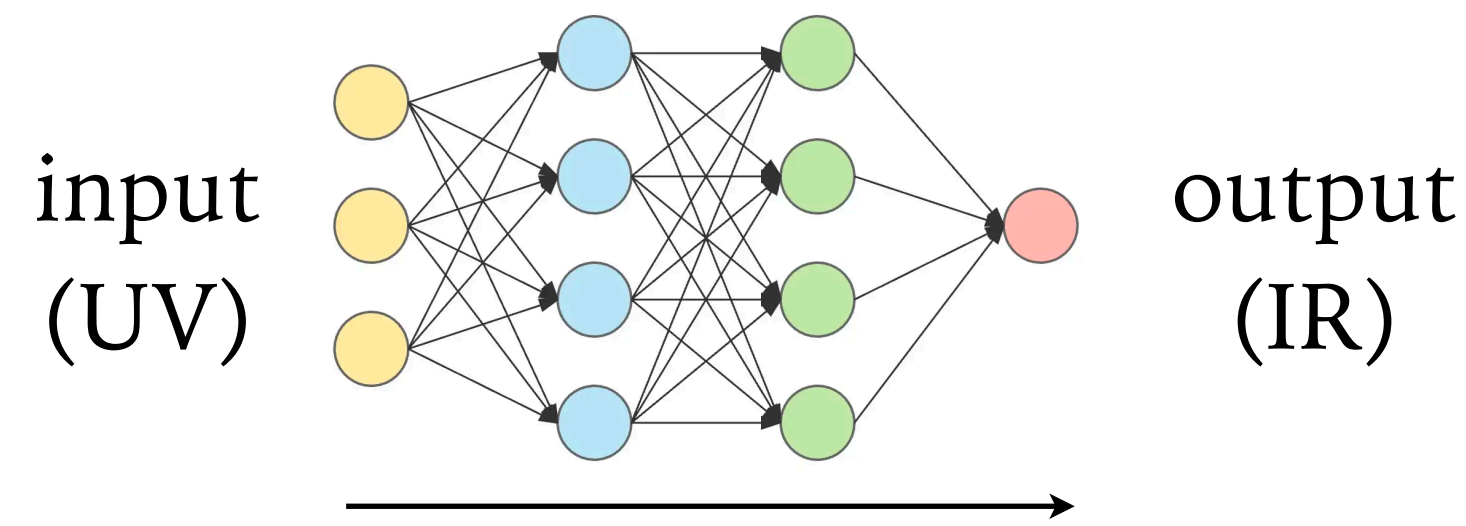
## 8-point



## 6-point

# Criticality

input
(UV)

output
(IR)

Exponential scaling (generic) $\leftrightarrow$ flow to trivial fixed point.

Tune to criticality $\Rightarrow$ power-law scaling $\leftrightarrow$ nontrivial fixed point.

Raghu et al '16. Poole et al '16. Schoenholz et al '16.

2-point correlator analysis:

$$\langle \mathcal{G}^{(\ell-1)}(\vec{x}_1, \vec{x}_2)\rangle \to \langle \mathcal{G}^{(\ell-1)}(\vec{x}_1, \vec{x}_2)\rangle + \delta\langle \mathcal{G}^{(\ell-1)}(\vec{x}_1, \vec{x}_2)\rangle$$

$$= \frac{C_W^{(\ell)}}{2}\left\langle \frac{\delta^2 \Delta(\vec{x}_1, \vec{x}_2)}{\delta\phi(\vec{y}_1)\delta\phi(\vec{y}_2)}\right\rangle_{\mathcal{K}_0^{(\ell-1)}} + \mathcal{O}\left(\frac{1}{n}\right)$$

$$\Rightarrow \ \delta\langle \mathcal{G}^{(\ell)}(\vec{x}_1, \vec{x}_2)\rangle = \ \bullet\!\!-\!\!\boxed{\delta}\!\!-\!\!\bullet \ = \ \sum_j \ \boxed{\begin{matrix}\delta \\ \phi_j \\ \Delta_j\end{matrix}} \ = \int d\vec{y}_1 d\vec{y}_2 \, \chi^{(\ell)}(\vec{x}_1, \vec{x}_2; \vec{y}_1, \vec{y}_2)\, \delta\langle \mathcal{G}^{(\ell-1)}(\vec{y}_1, \vec{y}_2)\rangle$$

susceptibility

Tune to criticality: $\qquad \chi^{(\ell)}(\vec{x}_1, \vec{x}_2; \vec{y}_1, \vec{y}_2)\Big|_{\mathcal{K}_0^{(\ell-1)}=\mathcal{K}^\star} = \frac{1}{2}\Big[\delta(\vec{x}_1 - \vec{y}_1)\,\delta(\vec{x}_2 - \vec{y}_2) + \delta(\vec{x}_1 - \vec{y}_2)\,\delta(\vec{x}_2 - \vec{y}_1)\Big]$

RG fixed point

# Structures of RG flow

Higher-point correlators?

Common structure:



$$\Rightarrow \quad \left(\frac{n_{\ell-2}}{n_{\ell-1}}\right)^{k-1} \frac{\delta V_{2k}^{(\ell)}(\vec{x}_1, \vec{x}_2; \ldots; \vec{x}_{2k-1}, \vec{x}_{2k})}{\delta V_{2k}^{(\ell-1)}(\vec{y}_1, \vec{y}_2; \ldots; \vec{y}_{2k-1}, \vec{y}_{2k})} = \text{sym.} \left[\prod_{k'=1}^{k} \chi^{(\ell)}(\vec{x}_{2k'-1}, \vec{x}_{2k'}; \vec{y}_{2k'-1}, \vec{y}_{2k'})\right]$$

same susceptibility that appeared in the 2-point correlator analysis!

Single criticality condition: $\chi^{(\ell)}(\vec{x}_1, \vec{x}_2; \vec{y}_1, \vec{y}_2)\big|_{\mathcal{K}_0^{(\ell-1)}=\mathcal{K}^\star} = \frac{1}{2}\left[\delta(\vec{x}_1 - \vec{y}_1)\,\delta(\vec{x}_2 - \vec{y}_2) + \delta(\vec{x}_1 - \vec{y}_2)\,\delta(\vec{x}_2 - \vec{y}_1)\right]$
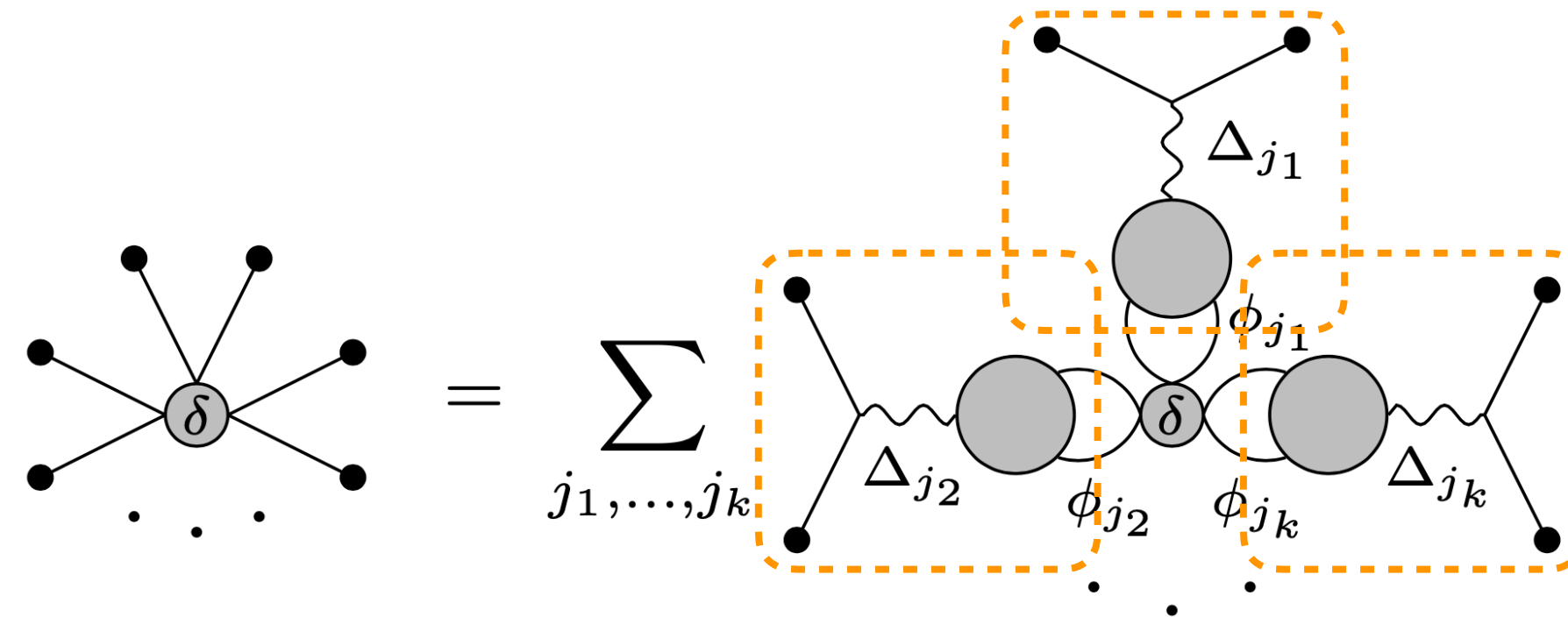
$\Rightarrow$ Power-law scaling for all connected correlators!

# Neural scaling laws & duality

Many ML models exhibit power law scaling of performance.

[Kaplan et al, 2001.08361] [Sharma, Kaplan, 2004.10802] [Bahri et al, 2102.06701]



approximately symmetric

Ising model of neural scaling laws? [Maloney, Roberts, Sully, 2210.16859]

# Large–N diagrammatics for neural scaling laws

$$\langle \text{Test loss} \rangle = \mathcal{R} \cdot (\mathcal{L}_1 + 2\mathcal{L}_2 + \mathcal{L}_3')$$



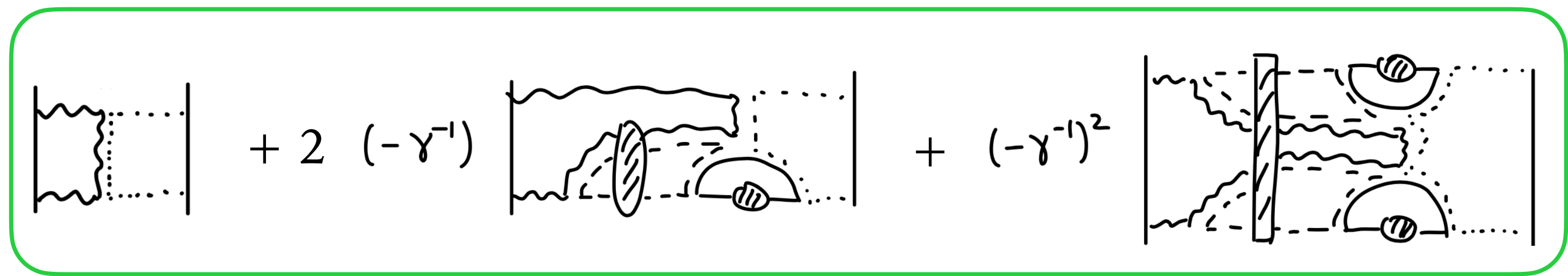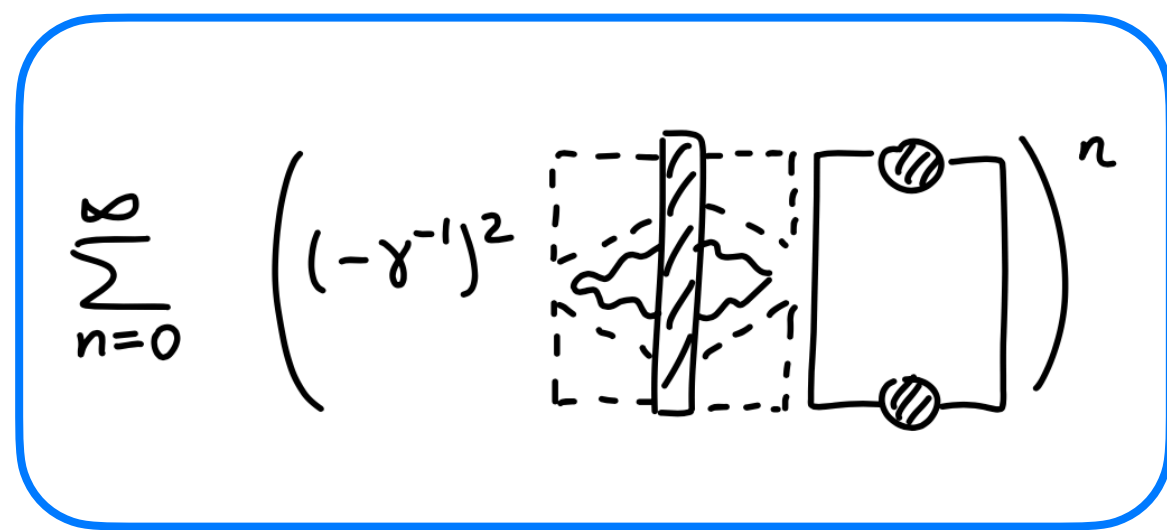Sets of diagrams are related by a duality transformation, e.g.

# Dreams

A theory of ~~everything~~ deep learning (opening the black box)?

Lee et al '17-19. Matthews et al '18. Yang '19-23.

Jacot, Gabriel, Hongler '18.

Antognini '19. Huang, Yau '19.

Yaida '19, '22. Hanin, Nica '19. Hanin '21, '22.

Dyer, Gur-Ari '19. Aitken, Gur-Ari '20. Andreassen, Dyer '20.

Naveh, Ringel et al '20, '21. Zavatone-Veth et al '21.

**Roberts, Yaida, Hanin '21. (Our work is largely inspired by this book.)**

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

# Dreams

A theory of ~~everything~~ deep learning (opening the black box)?

A new angle to learn about field theories?

$$\langle \phi(x_1) \dots \phi(x_k) \rangle = \int d\theta \, P(\theta) \, \phi_\theta(x_1) \dots \phi_\theta(x_k) = \frac{1}{Z} \int \mathcal{D}\phi \, e^{-S[\phi]} \, \phi(x_1) \dots \phi(x_k)$$

parameter/feature space
description
<span style="color:red">dual</span>
functional/sample space
description

Erbin, Lahoche, Samary '21, '22.

Bachtis, Aarts, Lucini '21.

Grosvenor, Jefferson (+ Erdmenger) '21.

**Halverson '21; + Maiti, Stoner '20, '21; + Demirtas, Schwartz '23.**



RABBIT

DUCK