

Space of Functions Computed by Deep-layered Machines

Bo Li, Alexander Mozeika and David Saad

Joint KIT/TUe Workshop on Neuromorphic High-Speed Communications (NeuCos)
9 December 2021

Outline

- Deep learning machines – successes and open questions
- Statistical mechanics of learning from examples and why entropy in function-space matters
- Related set-ups: Continuous/discrete weights, dense/sparse networks, correlated weights, convolutional neural networks, sensitivity to input perturbations, binarization/sparsification and finite-size effects
- Typical functions computed by DLM and recurrent networks
- Summary and future work

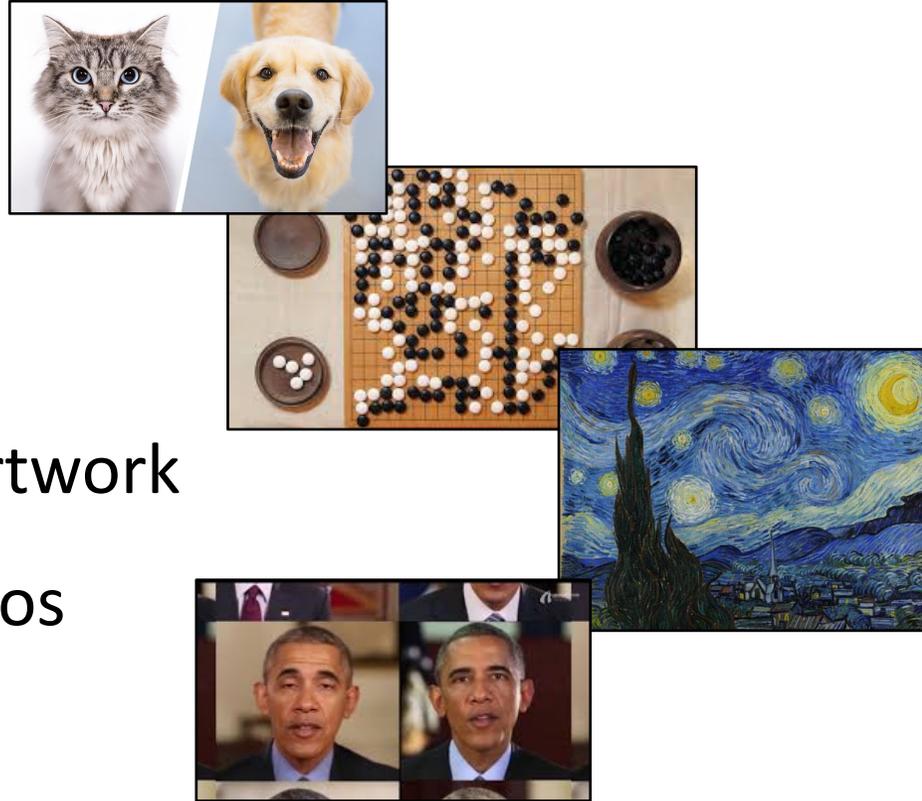
Bo Li and David Saad, Physical Review Letters **120**, 248301 (2018)

Bo Li and David Saad, Jour. Phys. A, **53**, 104002 (2020)

A. Mozeika, B. Li, and D. Saad, Physical Review Letters **125**, 168301, (2020)

Deep Learning Engineering Successes

- Computer vision
- Speech recognition
- Go, ATARI
- Composing music, artwork
- Generating fake videos



It is unclear:

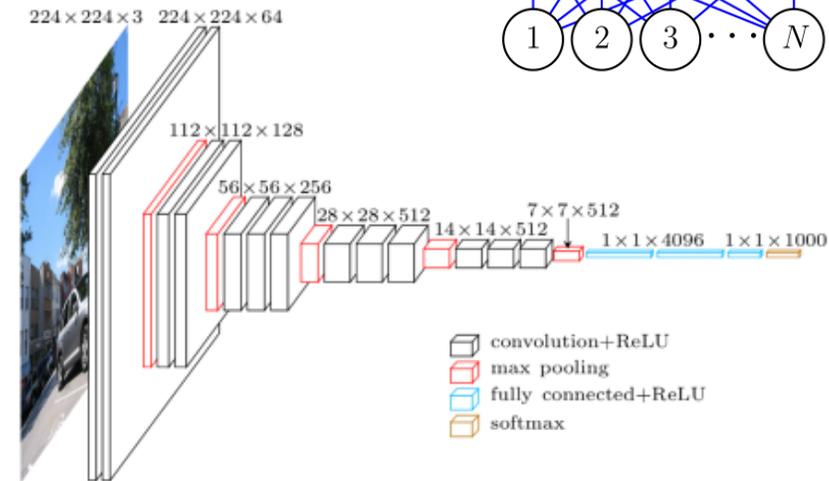
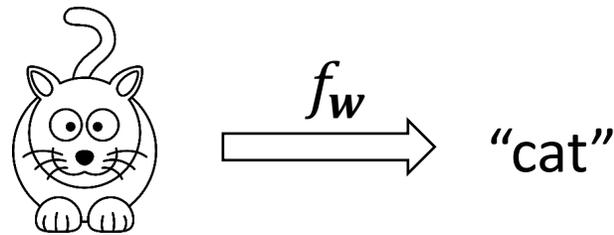
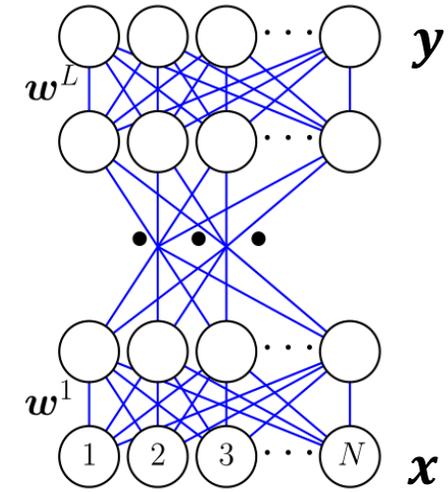
Which functions do they represent? How to optimize configurations? How to improve training? What do internal representations represent?

Deep Learning Machines

Implement an input-output mapping

$$y = f_w(x),$$

where the parameters w are to be estimated based on the training data $\{(\xi^\mu, \sigma^\mu)\}_{\mu=1,2,\dots,P}$ to perform a desired mapping.



We want to understand:

- (i) Their **generalization ability** even with numerous parameters
- (ii) Nature of the **internal representations**
- (iii) Which problems are easier/more difficult to solve?

Macroscopic Analysis – Typical Behavior

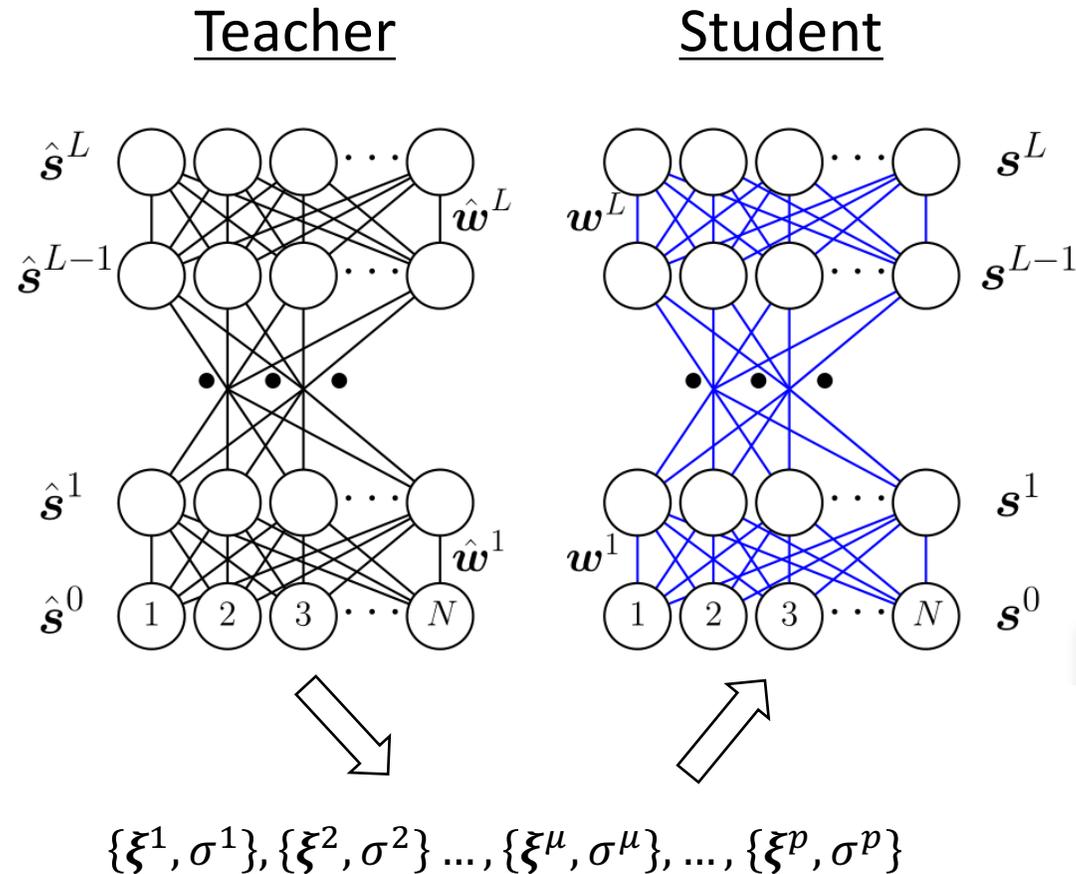
- Mapped to **disordered** systems of **infinite dimension**
- **Typical** behavior (in contrast to worse-case) of storage capacity and generalization curves (mostly single-layer)
- Technically quite involved (single or two-layer systems)
- Two-layer analysis – in the online setting only
- Input data structure and internal representations are rarely addressed

A. Engle and C. Van den Broeck, 2001,

D. Saad and S.A. Solla, Phys. Rev. Lett., 74, 4337, (1995); D. Saad and M. Rattray, Phys. Rev. Lett., 79, 2578 (1997)

M. Rattray, D. Saad and S. Amari, Phys. Rev. Lett., 81, 5461 (1998), S. Goldt, M. Mézard, F. Krzakala, and L. Zdeborová, Phys. Rev. X, 10, 041044 (2020)

Teacher-student Scenario for DLM?



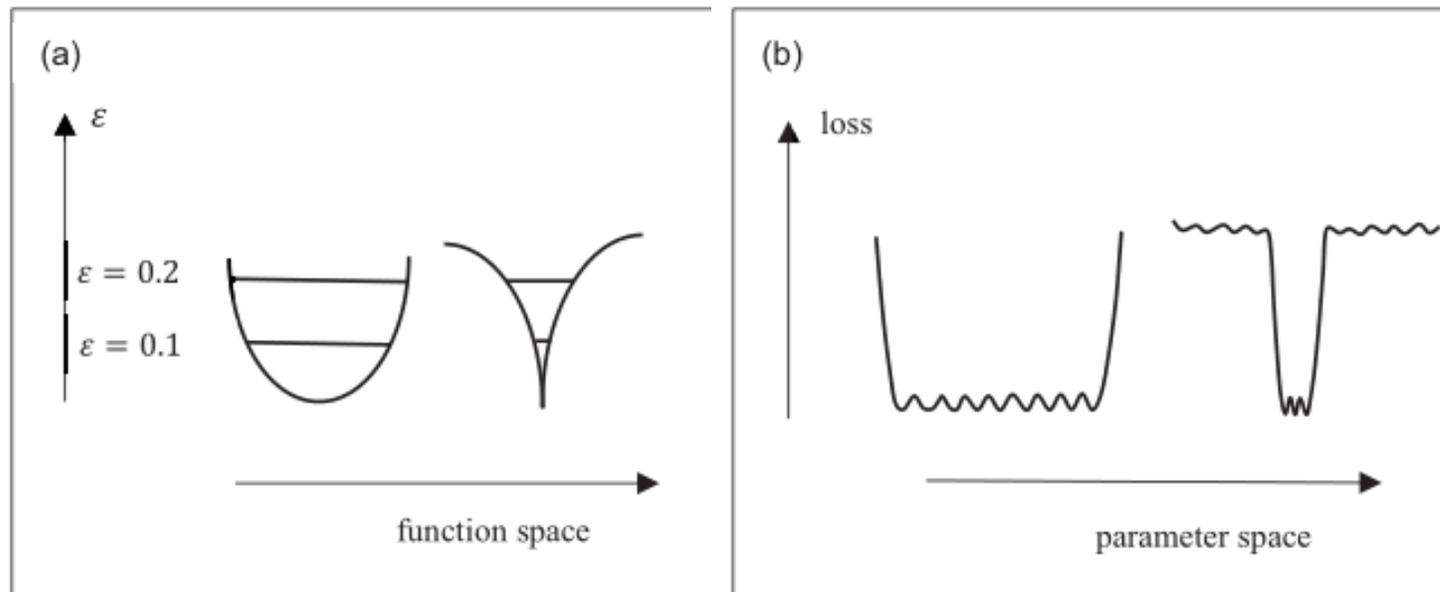
Difficulties:

- Constraints imposed by the examples (input-output pairs) on the hidden units are complex –**recursive nonlinear mapping**.
- Permutation, reflection and other **symmetries/invariances of hidden units**, no simple relation between teacher-student overlap and generalization error.

Function Space, Error and Entropy

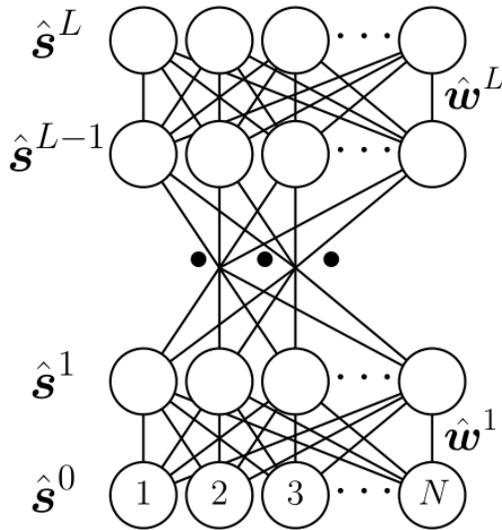
- We would like to approximate a reference/target function $f_{\hat{w}}$, as closely as possible from data.
- Given noisy data, sub-optimal training methods - more relevant to find *good approximations*. **How many such functions exist?**
- Conjecture - The entropy (log-volume) of functions at distance- ε away from $f_{\hat{w}}$ indicates how easy it is to obtain them.

Baldassi et al PRL 2015,
PNAS 2016, JSM 2020

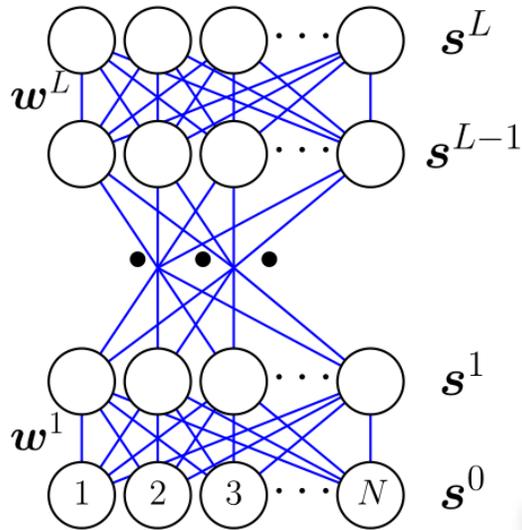


Exploring Function Space in DLM

Reference function $f_{\hat{w}}$

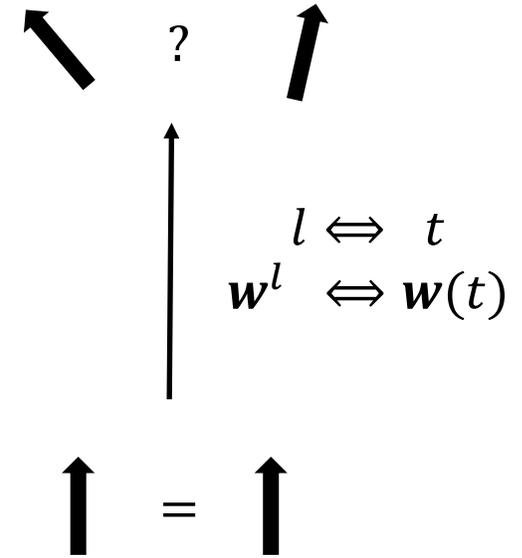


Perturbed function f_w



We map the DLM to **disordered spin systems** with discrete dynamics, $\hat{s}_i^l, s_i^l \in \{1, -1\}$, activation function is **sign function** $sgn(x)$.

The framework can be generalized to **real variables** and **other activation functions**.



Investigate the function sensitivity under small perturbations

$$w^l = \text{Perturb}(\hat{w}^l)$$

DLM as a Stochastic Dynamical System

- The layer evolution of two coupled DLMs:

$$P(\hat{\mathbf{s}}^l | \hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}, \beta) = \prod_i \frac{\exp \beta \hat{s}_i^l \hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1})}{2 \cosh \beta \hat{s}_i^l \hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1})}, P(\mathbf{s}^l | \mathbf{w}^l, \mathbf{s}^{l-1}, \beta) = \dots,$$

$\hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}) = \sum_j \hat{w}_{ij}^l \hat{s}_j^{l-1} / \sqrt{N}$, β is the inverse-temperature quantifying the noise level; deterministic rule in the zero-noise limit $\beta \rightarrow \infty$.

- Any observable is given by

$$\langle O \rangle := \sum_{\{\hat{\mathbf{s}}^l, \mathbf{s}^l\}} O \cdot P(\hat{\mathbf{s}}^0) \delta_{\hat{\mathbf{s}}^0, \mathbf{s}^0} \prod_l P(\hat{\mathbf{s}}^l | \hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}, \beta) \cdot P(\mathbf{s}^l | \mathbf{w}^l, \mathbf{s}^{l-1}, \beta),$$

summed over all the trajectories subject to the path measure.

- For discrete spins, the **overlap** between activities of the two systems is of interest

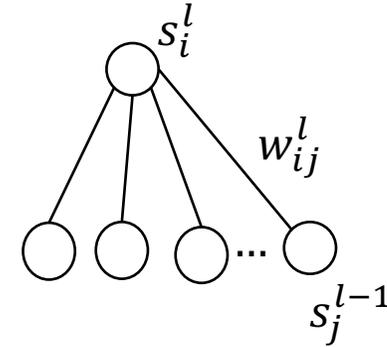
$$q^l(\hat{\mathbf{w}}, \mathbf{w}, \beta) = \frac{1}{N} \sum_i \langle \hat{s}_i^l s_i^l \rangle$$

Generating Functional Analysis

- Generating functional (characteristic function)

$$\Gamma[\hat{\psi}, \psi] := \left\langle \exp \left\{ -i \sum_{l,i} (\hat{\psi}_i^l \hat{s}_i^l + \psi_i^l s_i^l) \right\} \right\rangle,$$

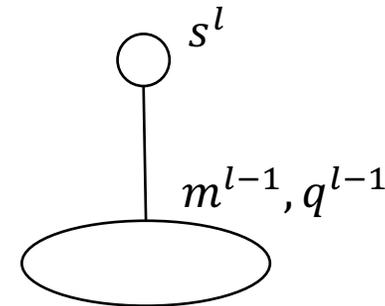
moments such as magnetization $\langle \hat{s}_i^l \rangle$ and overlap $\langle \hat{s}_i^l s_i^l \rangle$ can be obtained by differentiating $\Gamma[\hat{\psi}, \psi]$; angled brackets – average over all paths.



- Interested in the **typical** behavior of an ensemble of networks $\hat{\mathbf{w}} \sim P(\hat{\mathbf{w}})$, overbar – quenched average

$$\begin{aligned} \overline{\Gamma[\hat{\psi}, \psi]} &:= \sum_{\{\hat{\mathbf{w}}^l, \mathbf{w}^l\}} \Gamma[\hat{\psi}, \psi] P(\hat{\mathbf{w}}) P(\mathbf{w}) \\ &= \int \prod_l \frac{dQ^l dq^l}{2\pi/N} e^{N\Psi[\mathbf{q}, \mathbf{Q}]} \approx e^{N\Psi[\mathbf{q}^*, \mathbf{Q}^*]}, \text{ in the limit } N \rightarrow \infty, \end{aligned}$$

Represented by macroscopic order parameters; the saddle point $\mathbf{q}^*, \mathbf{Q}^* = \text{extr}_{\mathbf{q}, \mathbf{Q}} \Psi(\mathbf{q}, \mathbf{Q})$ satisfies certain self-consistent **mean-field** equation.

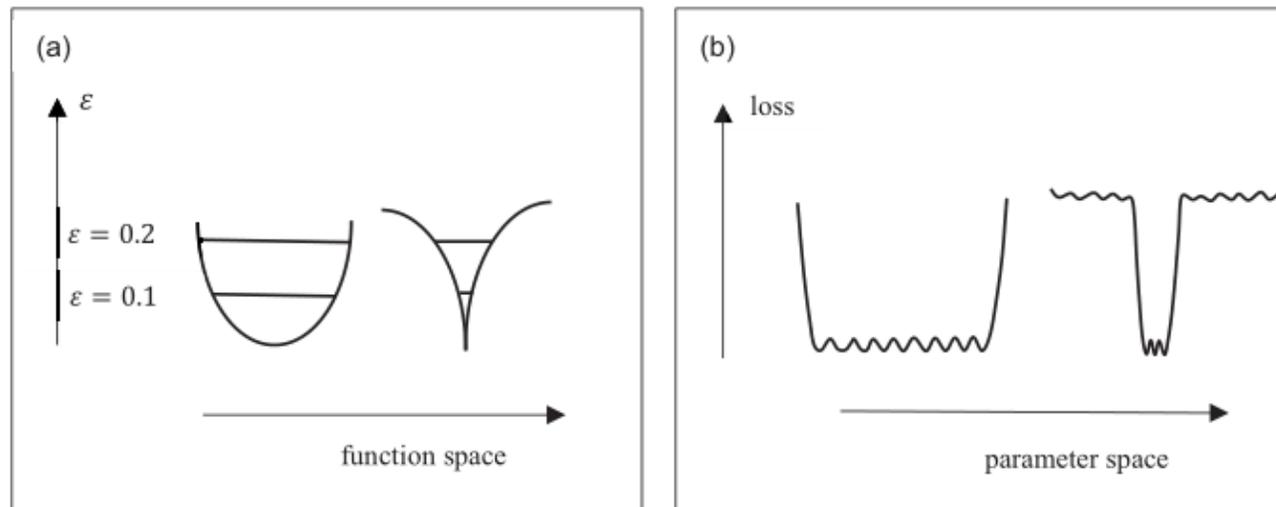


Function Error and Entropy

- For discrete spins, the **overlap** between internal representations the two systems is of interest $q^l(\hat{\mathbf{w}}, \mathbf{w}, \beta) = \frac{1}{N} \sum_i \langle \hat{s}_i^l s_i^l \rangle$, calculated using Generating Functional Analysis
- Function error is defined as the **expected Hamming distance** of output layers between $f_{\hat{\mathbf{w}}}$ and $f_{\mathbf{w}}$

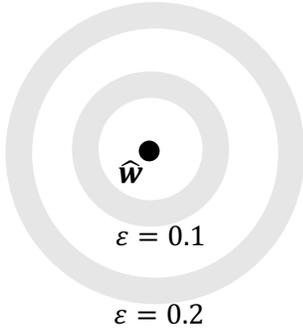
$$\varepsilon := \frac{1}{2N} \sum_{i=1}^N \overline{|\hat{s}_i^L - s_i^L|} = \frac{1}{2} (1 - q^L),$$

which provides a distance measure between $f_{\hat{\mathbf{w}}}$ and $f_{\mathbf{w}}$.



Fully-connected Networks – Continuous/Binary Weights

Consider fully-connected networks, with $P(\widehat{\mathbf{w}}_i^l) = \prod_j P(\widehat{w}_{ij}^l)$



Continuous weights:

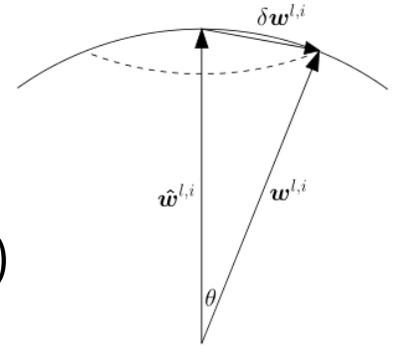
$$\widehat{w}_{ij}^l \sim \mathcal{N}(0, \sigma^2) \quad \text{Perturbation strength at layer } l$$

$$w_{ij}^l = \sqrt{1 - (\eta^l)^2} \widehat{w}_{ij}^l + \eta^l \delta w_{ij}^l$$

Binary weights:

$$P(\widehat{w}_{ij}^l) = \frac{1}{2} \delta(\widehat{w}_{ij}^l - 1) + \frac{1}{2} \delta(\widehat{w}_{ij}^l + 1)$$

$$P(w_{ij}^l) = (1 - p^l) \delta(w_{ij}^l - \widehat{w}_{ij}^l) + p^l \delta(w_{ij}^l + \widehat{w}_{ij}^l)$$



Typical overlaps: $q^l = \frac{2}{\pi} \sin^{-1}(\sqrt{1 - (\eta^l)^2} q^{l-1})$

$$q^l = \frac{2}{\pi} \sin^{-1}((1 - 2p^l) q^{l-1})$$

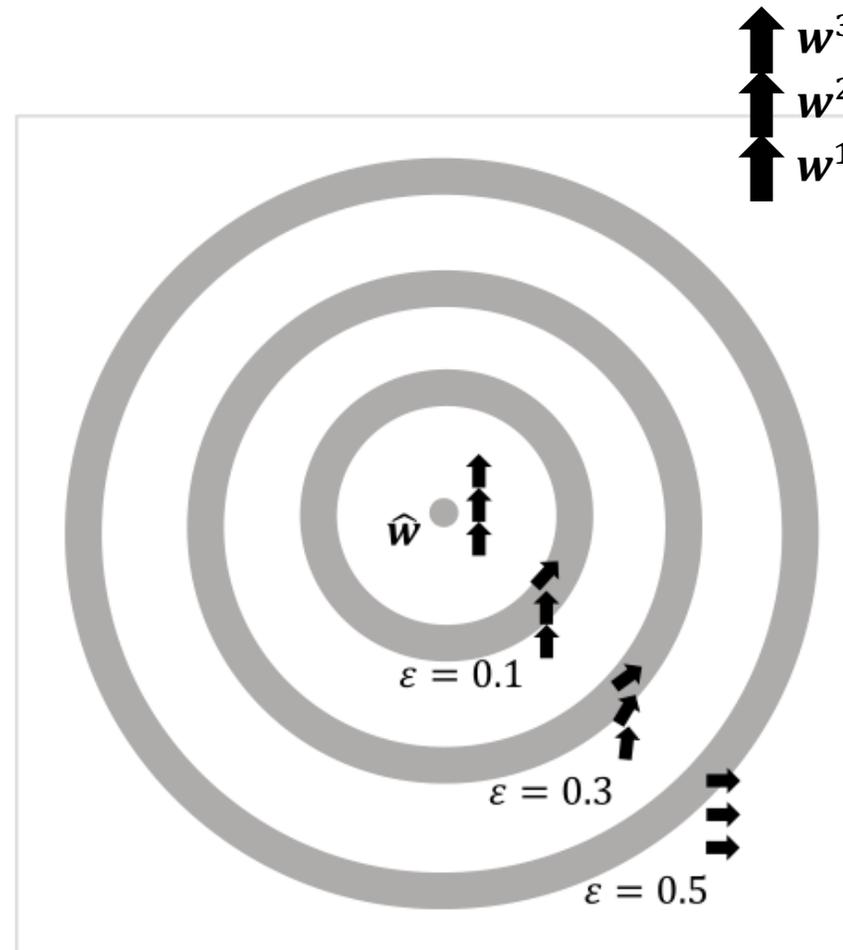
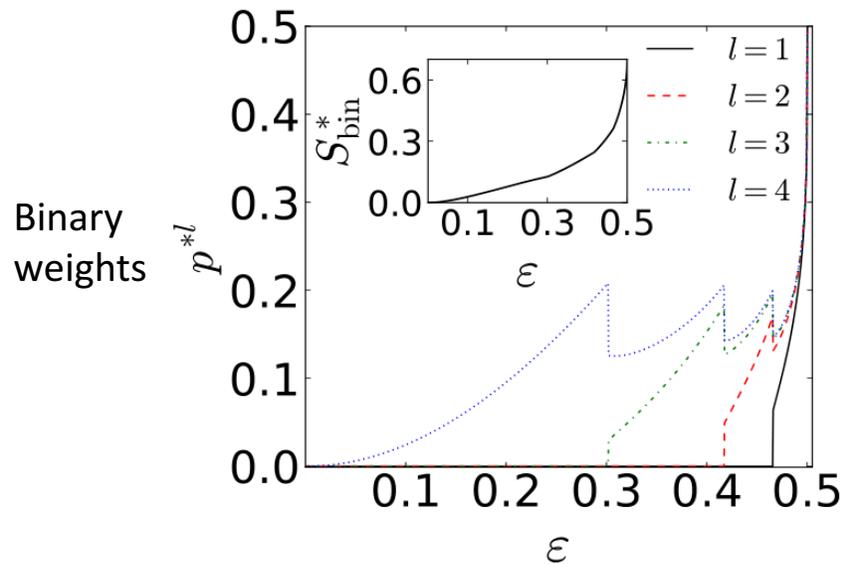
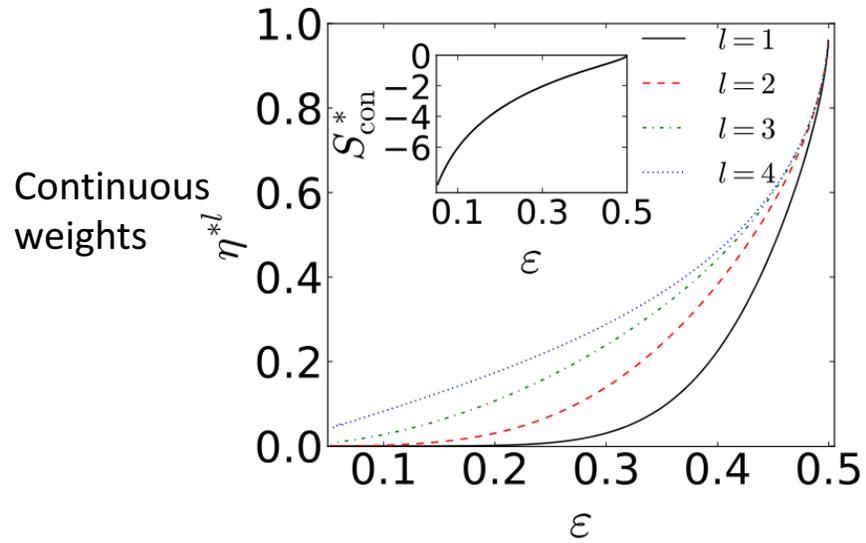
Entropy of f_w : $S_{\text{con}} = \frac{1}{L} \sum_l \log \eta^l$

$$S_{\text{bin}} = \frac{1}{L} \sum_l -p^l \log p^l - (1 - p^l) \log(1 - p^l)$$

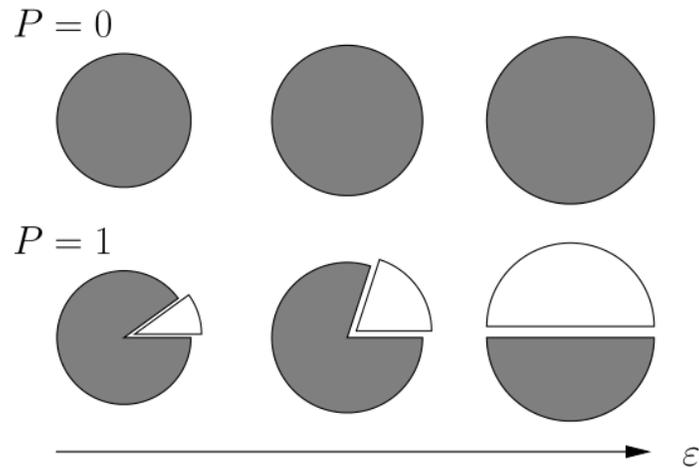
The distance- ε surface of f_w with volume $\Omega(\{\eta^l\}) = \exp N_p S_{\text{con}}(\{\eta^l\})$, is **exponentially** dominated by the maximum-entropy solutions when $N_p \rightarrow \infty$:

$$\eta^{*l} = \arg \max_{\eta^l} S_{\text{con}}(\{\eta^l\}), \quad \text{s.t. } q^L(\{\eta^l\}) = 1 - 2\varepsilon \quad 12$$

Earlier Layers Converge First When Decreasing ϵ



Approximate Generalization Curve (dense DLM with continuous weights)

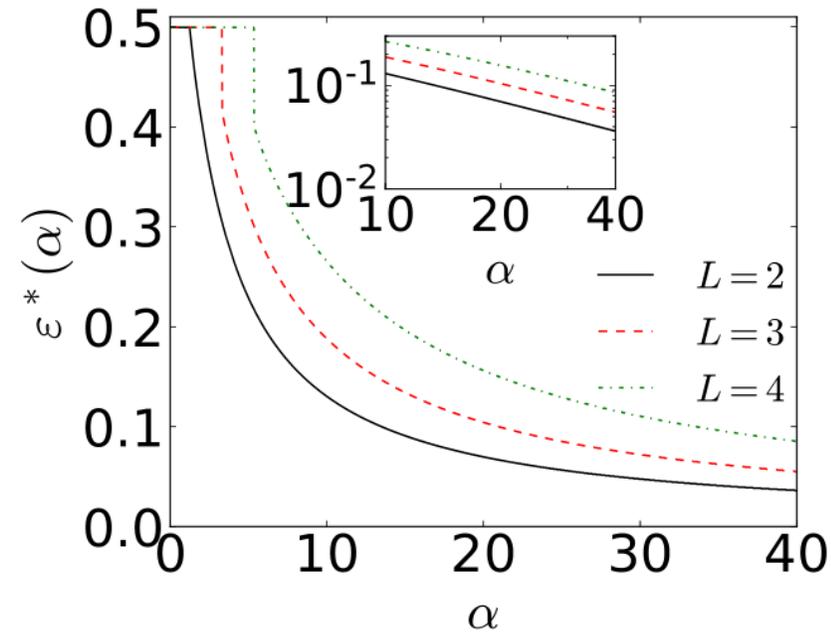


$$\Omega_\alpha(\varepsilon) = \Omega_{\text{tot}}(\{\eta^{*l}\}, \varepsilon) = \Omega_0(\varepsilon)(1 - \varepsilon)^P$$

$$P = \alpha LN^2 \quad \varepsilon^*(\alpha) = \operatorname{argmax}_\varepsilon \Omega_\alpha(\varepsilon)$$

Annealed theory of learning

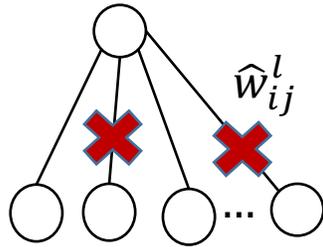
A. Engle and C. Van den Broeck, 2001



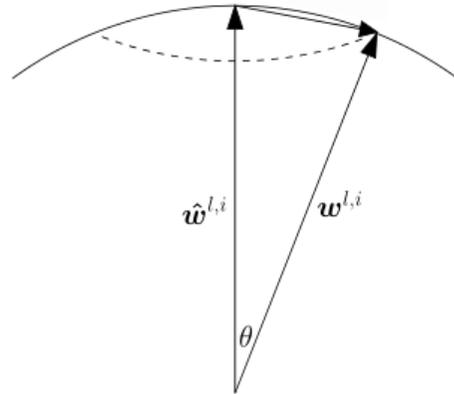
Relevant in small ε (large α) limit.

Perturbations Through Dilution/Discretization

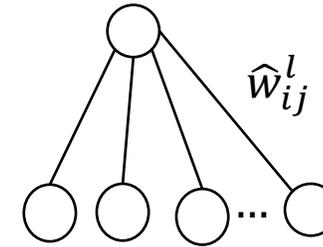
Weight disconnection



$$\theta^l = \sin^{-1} \sqrt{p^l}$$

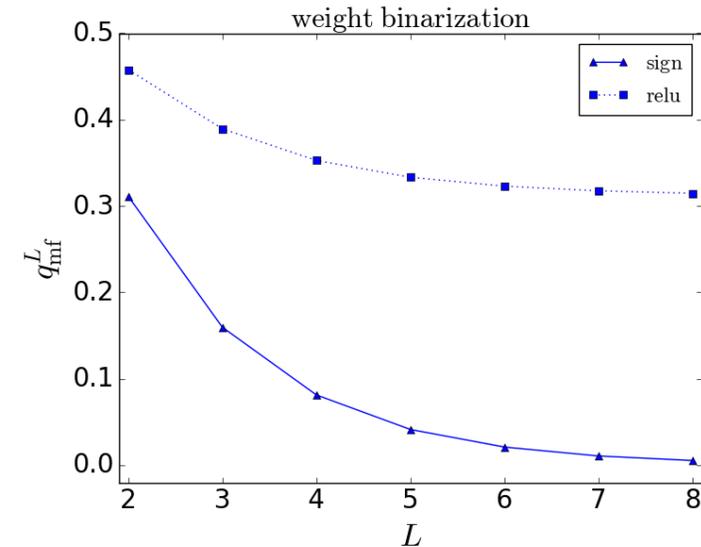
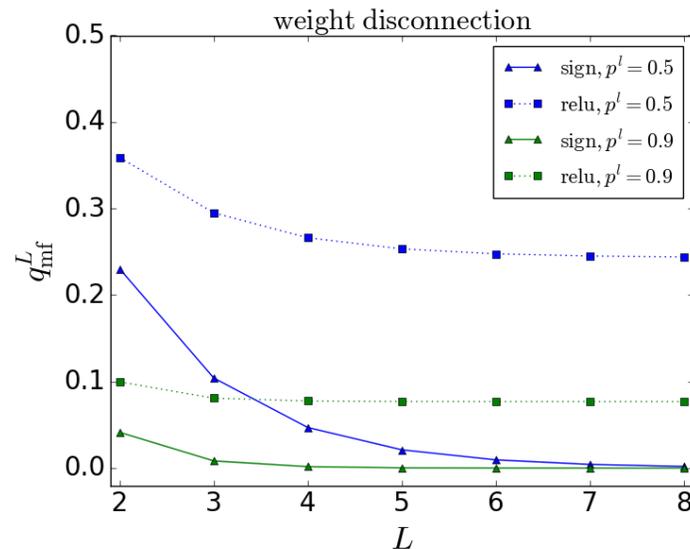


Weight discretization



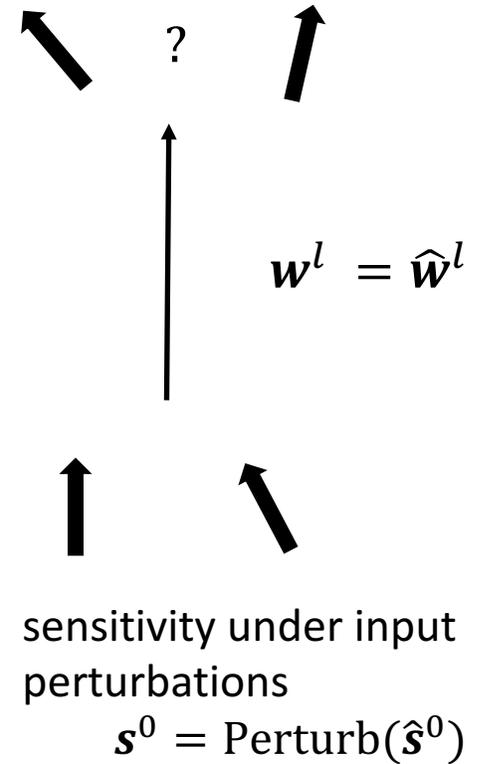
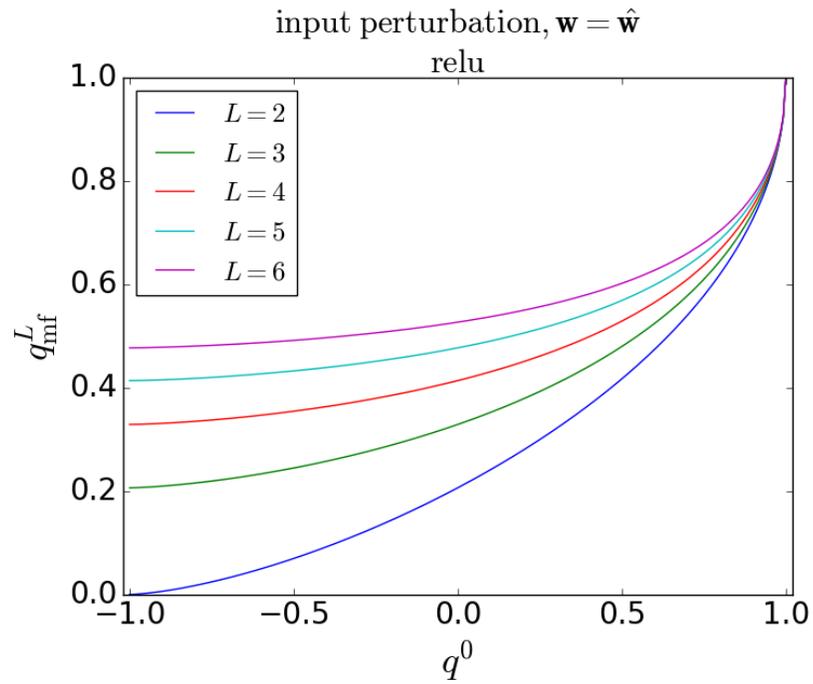
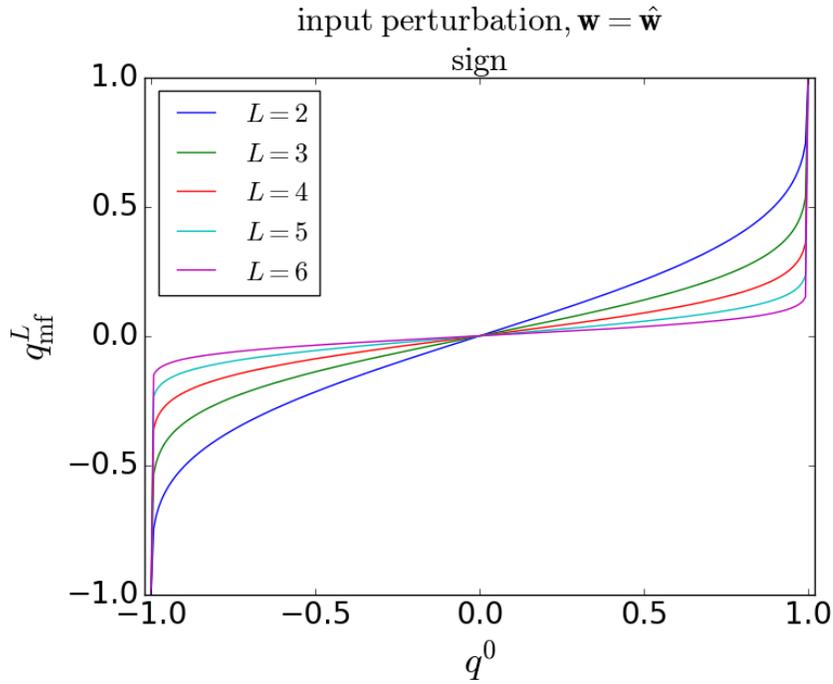
$$w_{ij}^l = \text{sgn}(\hat{w}_{ij}^l)$$

$$\theta^l = \cos^{-1} \sqrt{2/\pi} \approx 37^\circ \quad [\text{A. Anderson, et al., ICLR 2018}]$$



Deep ReLU networks with *random* weights are robust to disconnecting or binarizing weights.

Sign vs ReLU - Input Sensitivity

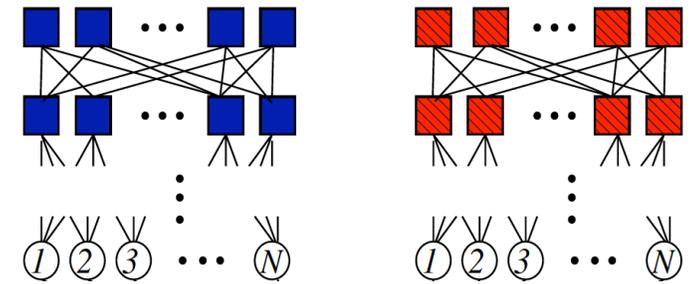


Deep ReLU networks with *random* weights compute simple functions.

[B. Poole et al., NIPS 2016]

Interim Summary – Other Setups we Investigated

- **Sparse architectures** - as before, where each node is **randomly** connected to k units in the previous layer and $\hat{w}_{ij}^l = 1$.



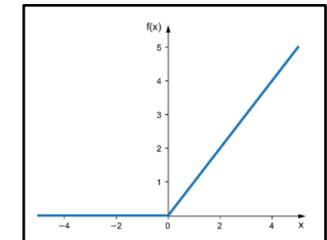
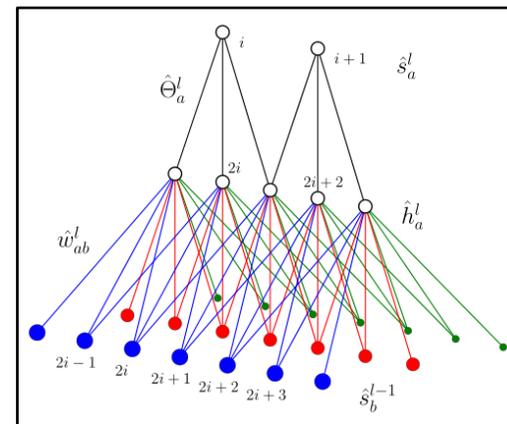
- **Weight dilution/discretization**– ReLU vs sign

- **Continuous variable values** and ReLU activation function $\phi(x) = \max(0, x)$.

- **Convolutional neural networks**

- **Correlated weights**

- **Input sensitivity** – ReLU vs sign



Space of Functions Generated by Sparse DLM

Consider Boolean functions computed by deep-layered machines,

$$f : \{+1, -1\}^n \rightarrow \{+1, -1\}.$$

E.g., for $n = 2$, there are $2^n = 4$ possible input patterns in total, which are

$$\begin{aligned}x_0 &= (+1, +1), & x_1 &= (+1, -1), \\x_2 &= (-1, +1), & x_3 &= (-1, -1).\end{aligned}$$

The Boolean function $f(\cdot)$ is fully specified by the string of length 4

$$\mathbf{f} \equiv (f(x_0), f(x_1), f(x_2), f(x_3)).$$

Boolean Input-Output Relation in DLM

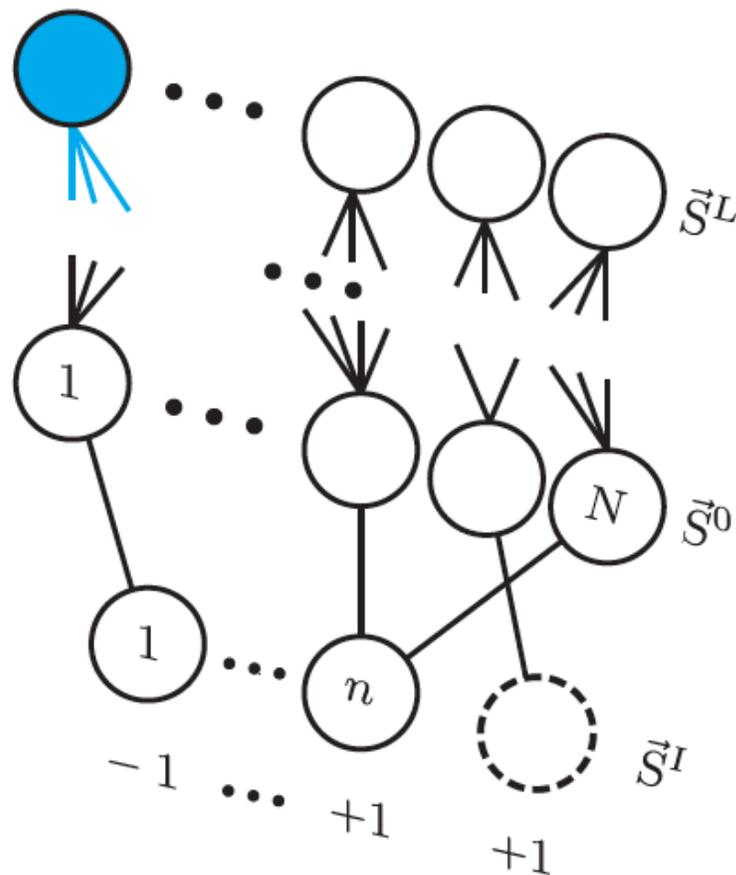
Computation of a specific input pattern $\vec{s} \in \{+1, -1\}^n$

$$P(\vec{S}^L | \vec{s}) = \sum_{\vec{S}^{L-1} \dots \vec{S}^0} P(\vec{S}^0 | \vec{s}) \prod_{\ell=1}^L P(\vec{S}^\ell | \vec{S}^{\ell-1}),$$

where

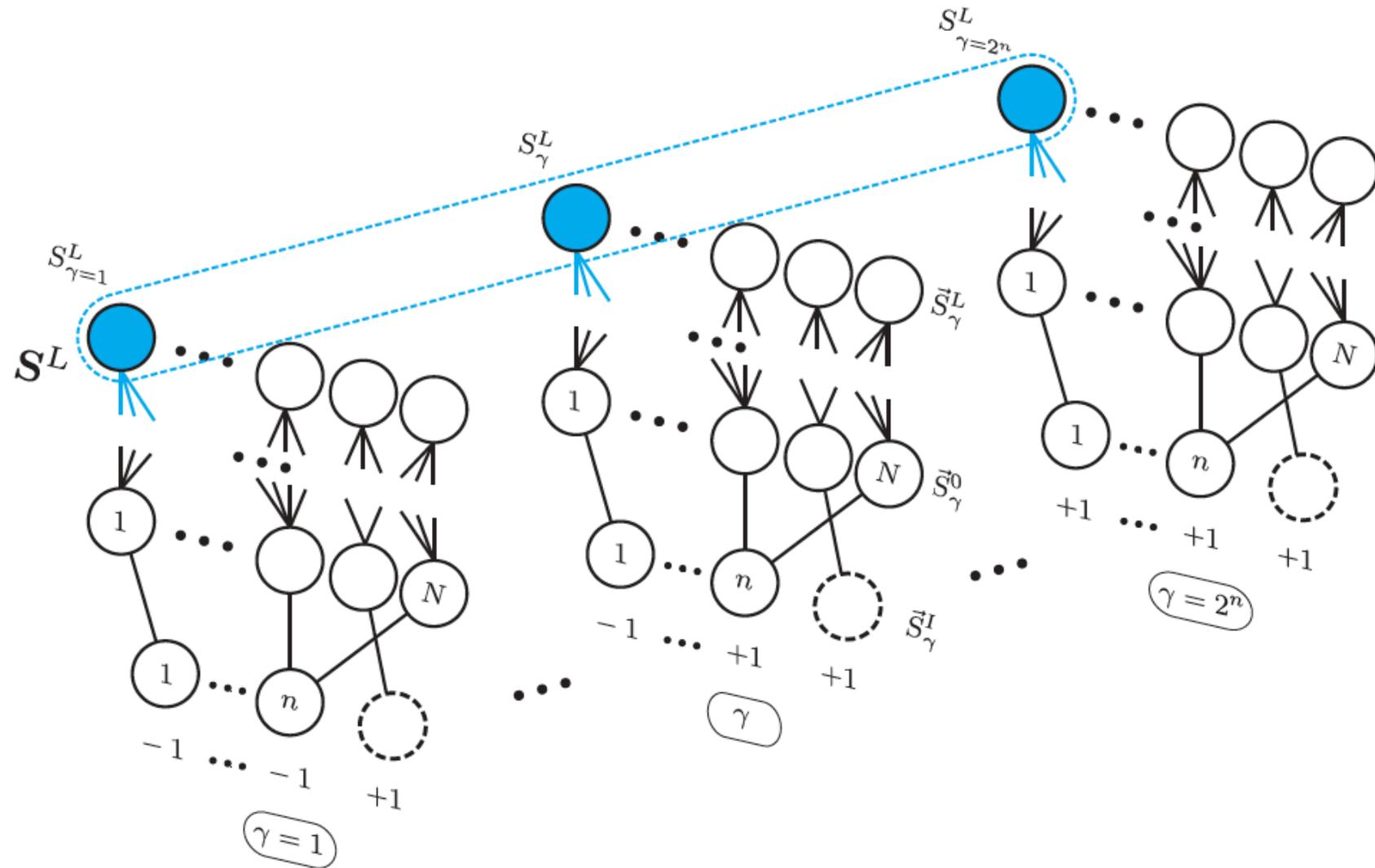
$$P^0(\vec{S}^0 | \vec{s}) = \prod_{i=1}^N \delta[S_i^0, S_{n_i}^I(\vec{s})],$$

$$P(S_i^\ell | \vec{S}^{\ell-1}) = \delta[S_i^\ell, \alpha_i^\ell(\vec{S}^{\ell-1})].$$



Functions Generated by Layered Networks

Computation of all possible input patterns: $\prod_{\gamma=1}^{2^n} P(\vec{S}_\gamma^L | \vec{s}_\gamma)$



The Framework – Probability of Functions

Distribution of functions computed on the final layer

$$P_N^L(\mathbf{f}) = \frac{1}{N} \sum_{i=1}^N \left\langle \prod_{\gamma=1}^{2^n} \delta(f_\gamma, S_{i,\gamma}^L) \right\rangle.$$

To compute macroscopic observables such as $P_N^L(\mathbf{f})$, introduce the generating functional

$$\overline{\Gamma[\{\psi_{i,\gamma}^\ell\}]} = \frac{1}{\sum_{\{S_{i,\gamma}^\ell\}_{\forall \ell, i, \gamma}}} \prod_{\gamma=1}^{2^n} P(\vec{S}_\gamma^L | \vec{s}_\gamma) e^{-i \sum_{\ell, i, \gamma} \psi_{i,\gamma}^\ell S_{i,\gamma}^\ell}.$$

GF Analysis – Saddle Point Equations

Weight disorder: $W_{ij}^\ell \sim \mathcal{N}(0, \sigma^2)$

By introducing the overlap $q_{\gamma\gamma'}^{\ell,\ell'} = (1/N) \sum_i \overline{\langle S_{i,\gamma}^\ell S_{i,\gamma'}^{\ell'} \rangle}$, the GF can be computed as

$$\bar{\Gamma} = \int \{d\mathbf{q} d\mathbf{Q}\} e^{N\Psi[\mathbf{q}, \mathbf{Q}]}.$$

The saddle points of $\Psi[\mathbf{q}, \mathbf{Q}]$ give rise to the typical behaviors

$$\frac{\partial \Psi}{\partial \mathbf{q}} = 0, \quad \frac{\partial \Psi}{\partial \mathbf{Q}} = 0.$$

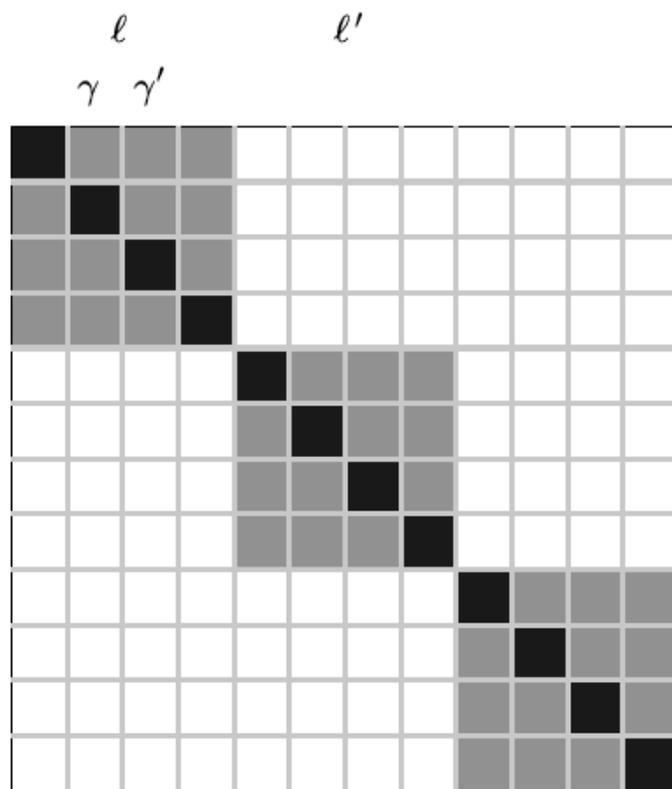
Layered vs Recurrent Networks

Overlaps for both architectures:

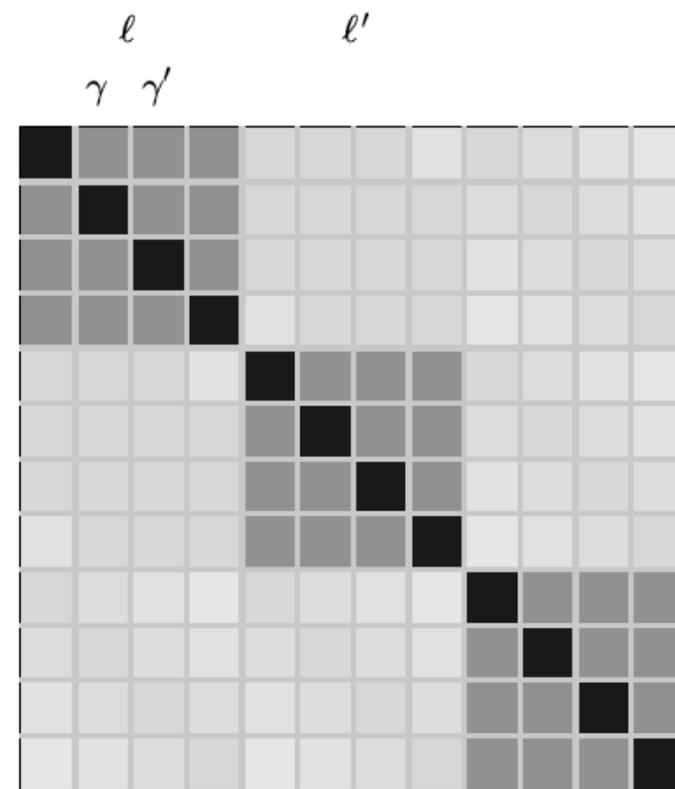
$$q_{\gamma\gamma'}^{l,l'} = \int d\mathbf{H} \alpha^l(h_\gamma^l) \alpha^{l'}(h_{\gamma'}^{l'}) \mathcal{N}(\mathbf{H}|\mathbf{0}, \mathbf{C}), \quad l, l' > 0.$$

Covariance matrix $\mathbf{C} \propto \mathbf{q}$:

The elements admit the same values for **both** architectures when $l = l'$.



layer-dependent



recurrent

Layered vs Recurrent - Functions

To examine the distribution of functions, we only need to consider single-layer activities! (Remind: $P_N^L(\mathbf{f}) = \frac{1}{N} \sum_i \langle \prod_{\gamma=1} \delta(f_\gamma, S_{i,\gamma}^L) \rangle$.)

- For **both** architectures of neural networks,

$$P^L(\mathbf{f}) = \int d\mathbf{h} \mathcal{N}(\mathbf{h}|\mathbf{0}, \mathbf{c}^L) \prod_{\gamma=1}^{2^n} \delta[f_\gamma, \text{sgn}(h_\gamma)].$$

- For **both** architectures of Boolean circuits,

$$P^{\ell+1}(\mathbf{f}) = \sum_{\mathbf{f}_1, \dots, \mathbf{f}_k} \left\{ \prod_{j=1}^k P^\ell(\mathbf{f}_j) \right\} \prod_{\gamma=1}^{2^n} \delta[f_\gamma, \alpha(f_{1,\gamma}, \dots, f_{k,\gamma})],$$

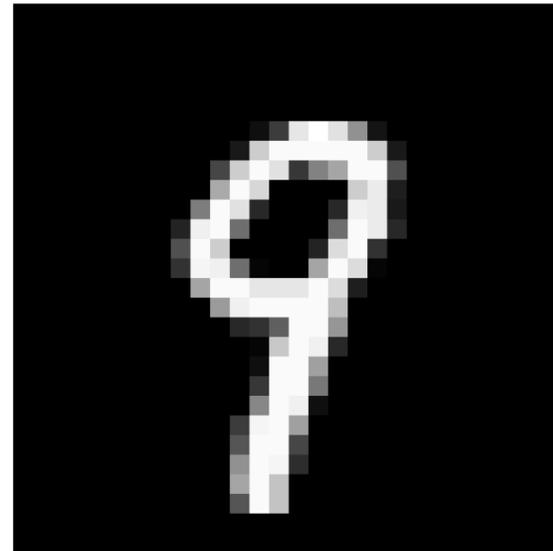
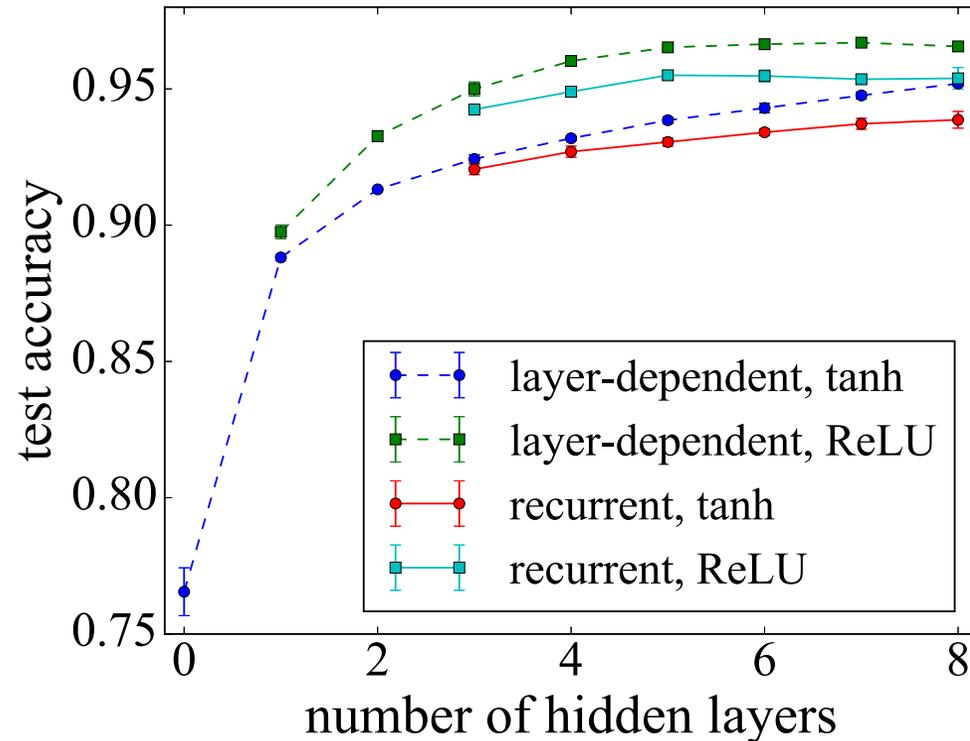
Numerical Results

Computation with weight-sharing across layers to save parameters?

Though the theory only applies to **random** machines, we consider training experiments with neural networks on MNIST data.

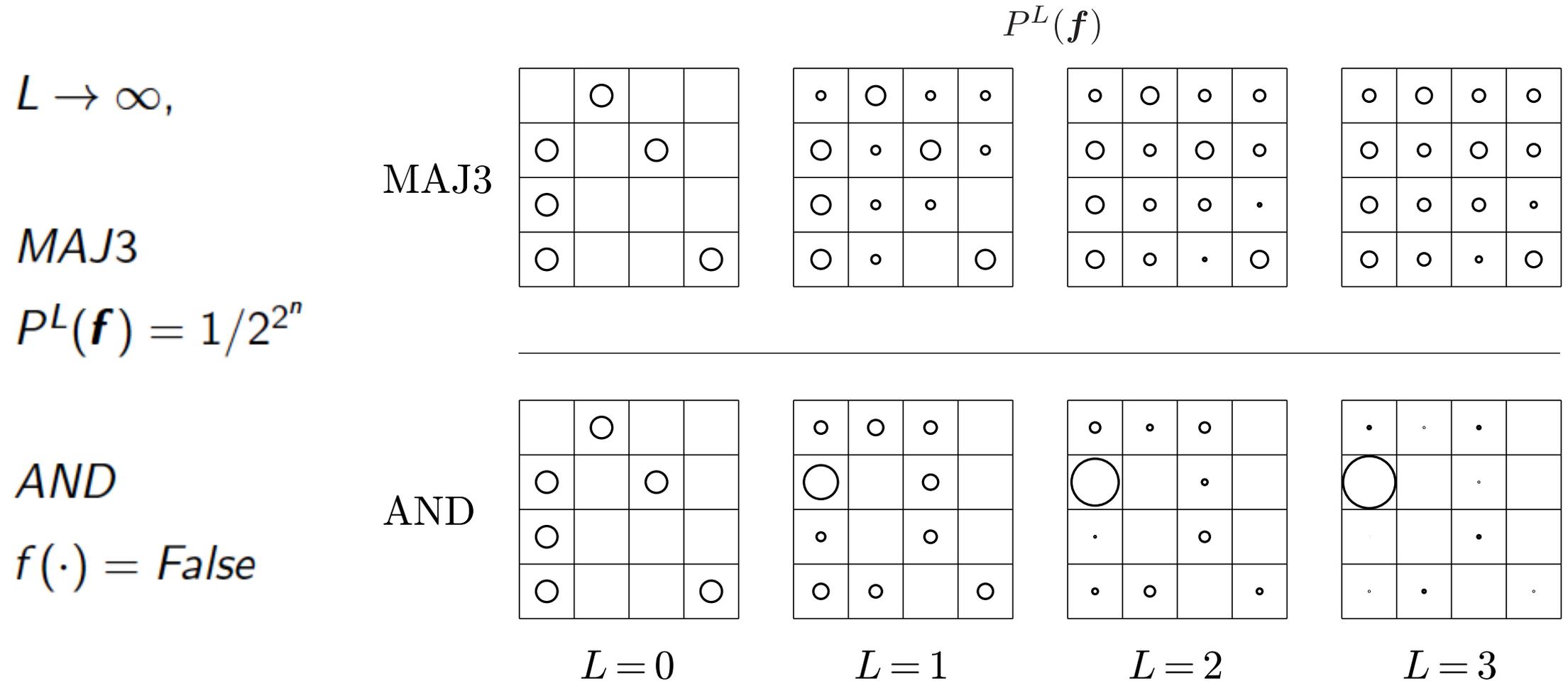
For recurrent architectures

$$\vec{S}^0 \xrightarrow{W^{\text{in}}} \vec{S}^1 \xrightarrow{W^{\text{hid}}} \vec{S}^2 \xrightarrow{W^{\text{hid}}} \dots \xrightarrow{W^{\text{hid}}} \vec{S}^{L-1} \xrightarrow{W^{\text{out}}} \vec{S}^L$$



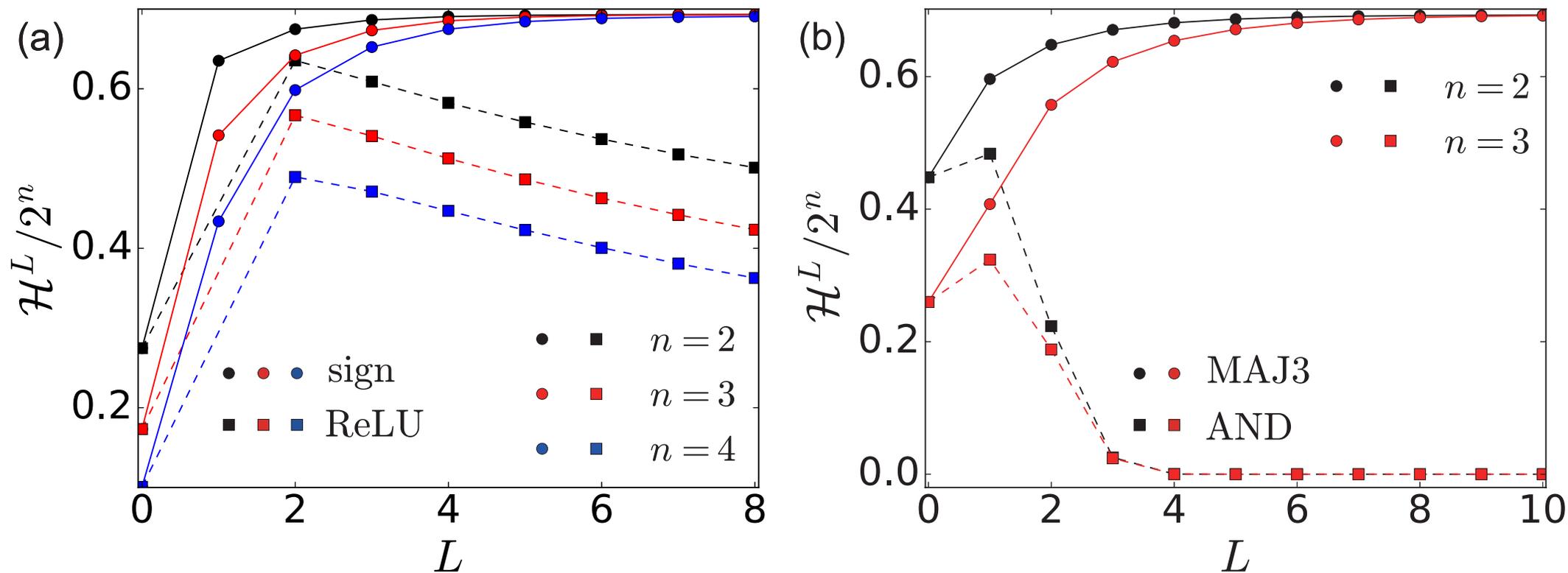
Entropy In Sparse Layered Networks

Effect of gate $\alpha(\cdot)$. Consider $\vec{S}^l = (\vec{s}, -\vec{s}, 1, -1)$.



Entropy of Functions – ReLU vs sign

Shannon entropy of Boolean functions $\mathcal{H}^L = -\sum_{\mathbf{f}} P^L(\mathbf{f}) \log P^L(\mathbf{f})$



Simplicity bias of neural networks with ReLU activation when initialized randomly, which arguably plays a role in their generalization ability [G. Valle-Perez et al., ICML 2019; G. De Palma et al., NeurIPS 2019].

Current and Future Work

- We proposed a framework for analyzing functions computed by random DLM.
- Show equivalence between the space of functions generated by random DLM and recurrent architectures.
- Possible computation with a less parameters by weight/connection sharing (may sacrifice accuracy)
- Depending on gate/activation functions, simpler or more complex functions are being represented as layer depth increases.

Future work:

- The role of **over-parametrization** in function landscape, error and generalization.
- Exploring the effect of non-trivially **correlated inputs**, e.g. from a generative model, and go beyond the random reference functions.
- Optimizing **variable hidden layer size**.