

Replica theory on a deep neural network

Cybermedia Center & Department of Physics,
Osaka Univ.

Hajime Yoshino

H.Yoshino,
arXiv1910.09918



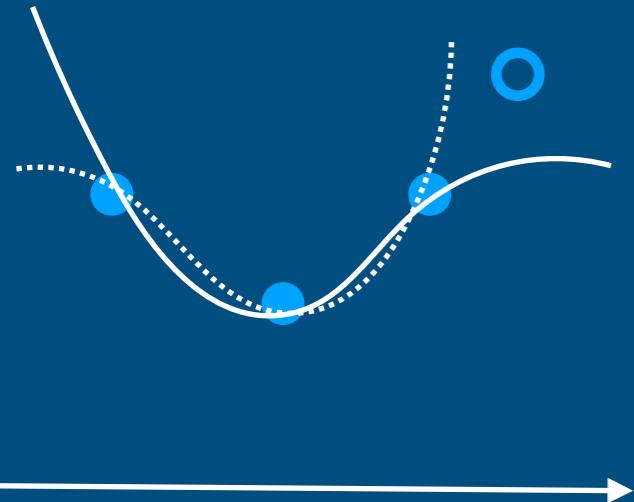
Deep networks are typically over-parametrized.

e.g.

data size	# of parameters
10^6	$\ll 10^8$

G. Carleo, et. al, arXiv:1903.10563v1

over fitting

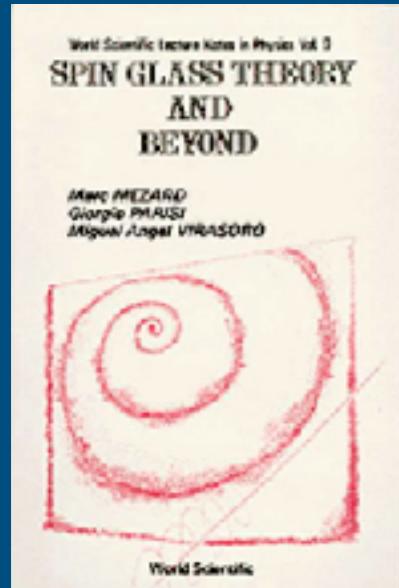


Common sense

1. poor generalization
2. many local minimum
learning is difficult
(glassy dynamics)

Empirical observations
on deep networks

1. generalization is not bad
???
2. learning is not too difficult
???



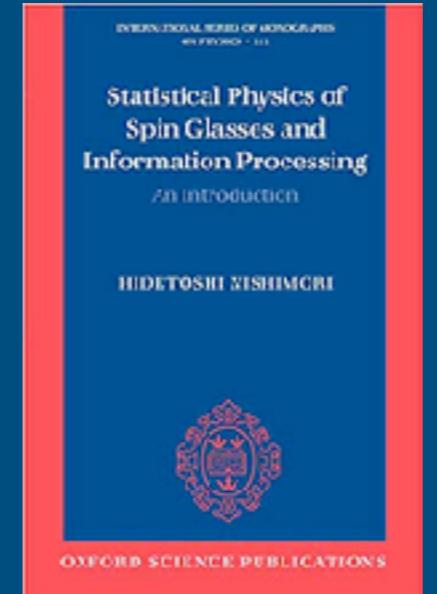
1987



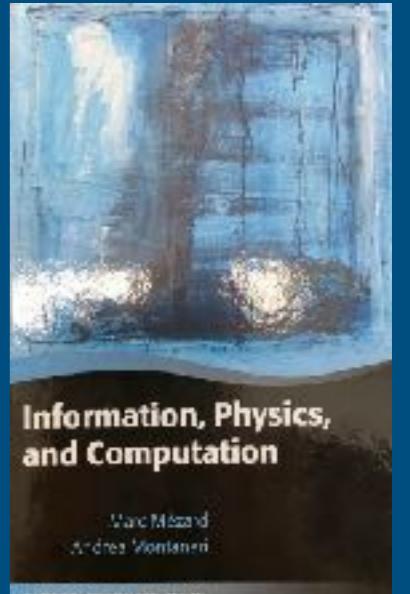
1991



1999



2001



2008

From spin glass to structural glass

p-spin to RFOT

more on replicas

$d \rightarrow \infty$

shear

Kirkpatrick-Thirumalai Wolynes (1989)

Franz-Parisi (1995), Monasson (1995), Mezard-Parisi (1999)

Parisi-Zamponi (2010), Charbonneau-Kurchan-Parisi-Urbani-Zamponi (2014)

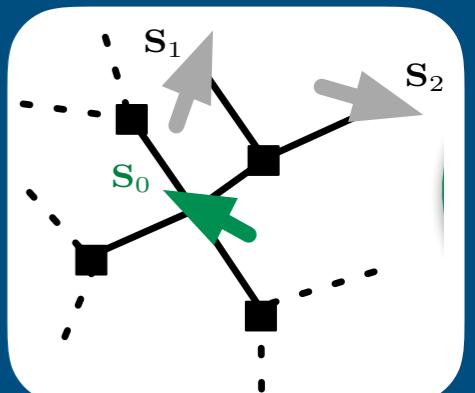
Yoshino-Mezard (2010), Yoshino-Zamponi (2014), Rainone-Urbani-Yoshino-Zamponi (2015)

Jin-Yoshino (2017), Jin-Urbani-Zamponi-Yoshino (2018)

... and back to p-spins... but with $M \rightarrow \infty$ spin components and without quenched disorder

Using explicit RSB : Parisi-Virasoro (1989)

H. Yoshino, SciPost Phys. 4 (6), 040 (2018)

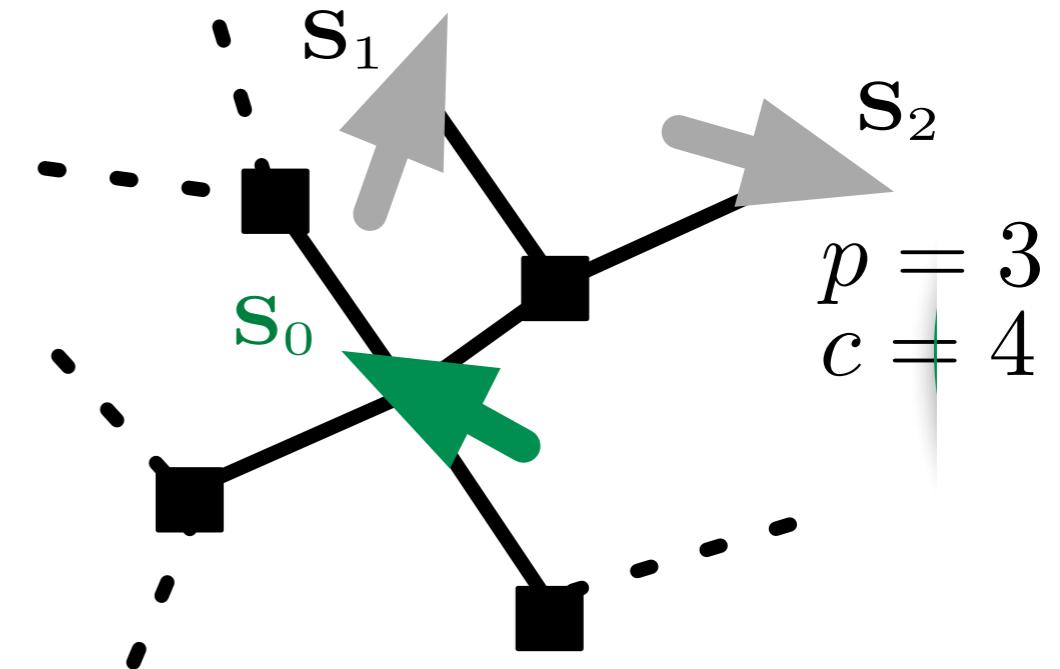


“Disorder-free” vector p-spin models on tree

H.Yoshino, SciPost Phys. 4 (6), 040 (2018)

$$\mathbf{S}_i = (S_i^1, S_i^2, \dots, S_i^M) \\ i = 1, 2, \dots, N \quad |S_i|^2 = M$$

continuous or Ising $S_i^\mu = \pm 1$



Hamiltonian

$$H = - \sum_{\blacksquare} V(r_{\blacksquare}).$$

Locally tree like lattice
with connectivity c

$$c = \alpha M$$

“gap” $r_{\blacksquare} = \delta - \frac{1}{\sqrt{M}} \sum_{\mu=1}^M S_{1(\blacksquare)}^\mu S_{2(\blacksquare)}^\mu \cdots S_{p(\blacksquare)}^\mu$

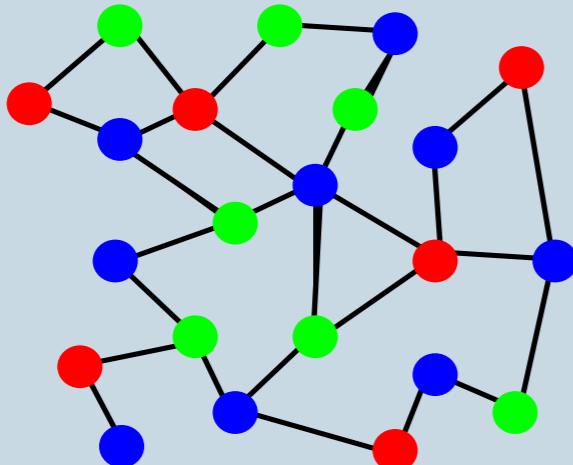
of factor nodes

$$N_{\blacksquare} = Nc/p = NM\alpha/p < N^p/p!$$

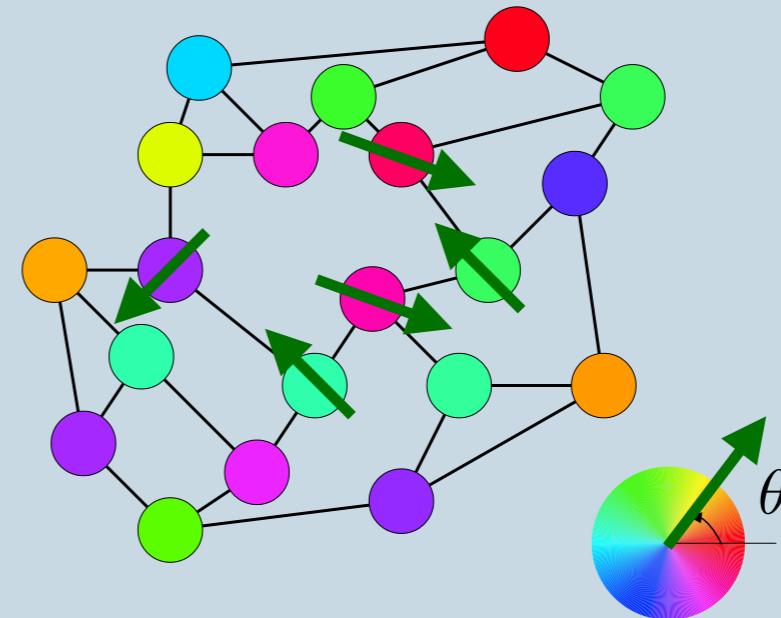
“Inter-mediate sparseness “: high connectivity but not “global coupling”

“Vectorial” constraint satisfaction problems

standard discrete coloring



“continuous” version



antiferromagnetic Potts model

$$H = \sum_{i,j} \delta_{q_i, q_j}$$

“repulsive” vectorial spin model

$$H = \sum_{i,j} V \left(\delta - \frac{\mathbf{S}_i \cdot \mathbf{S}_j}{\sqrt{M}} \right)$$

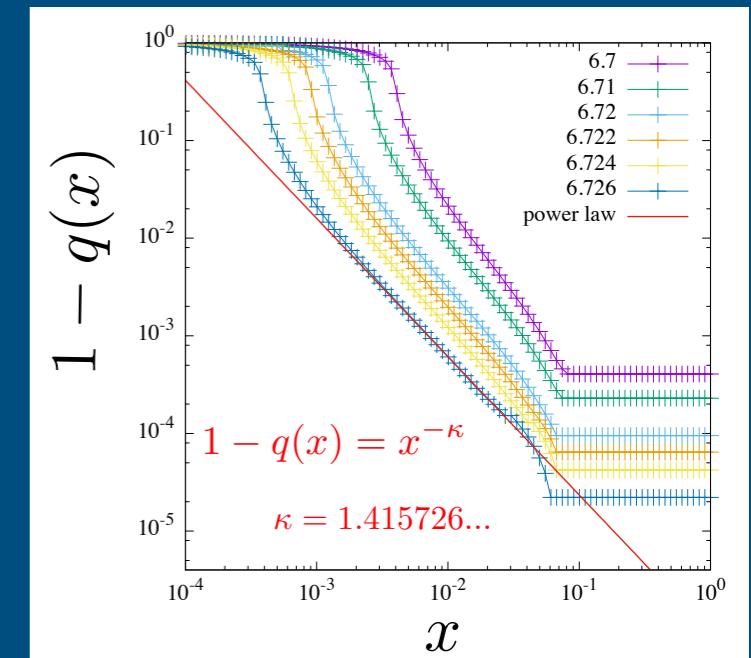
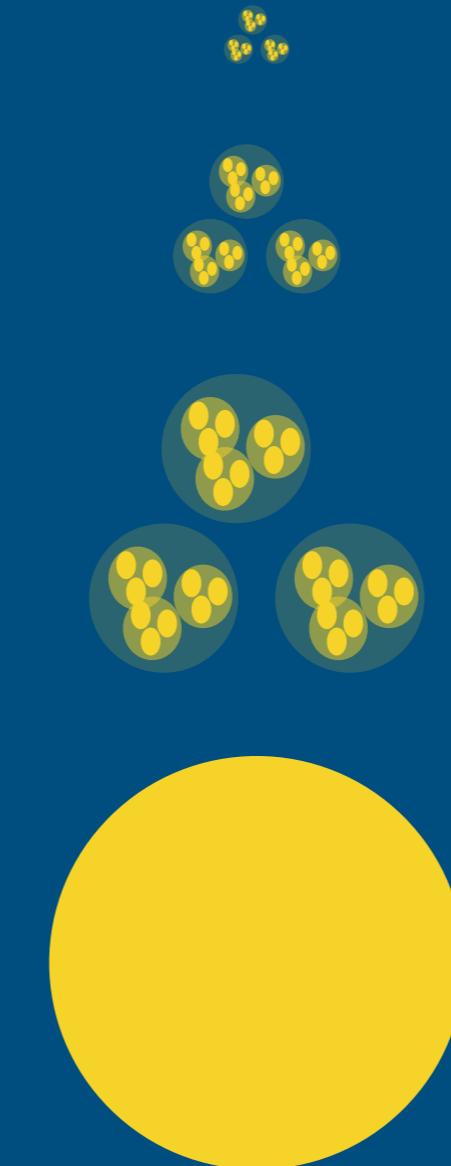
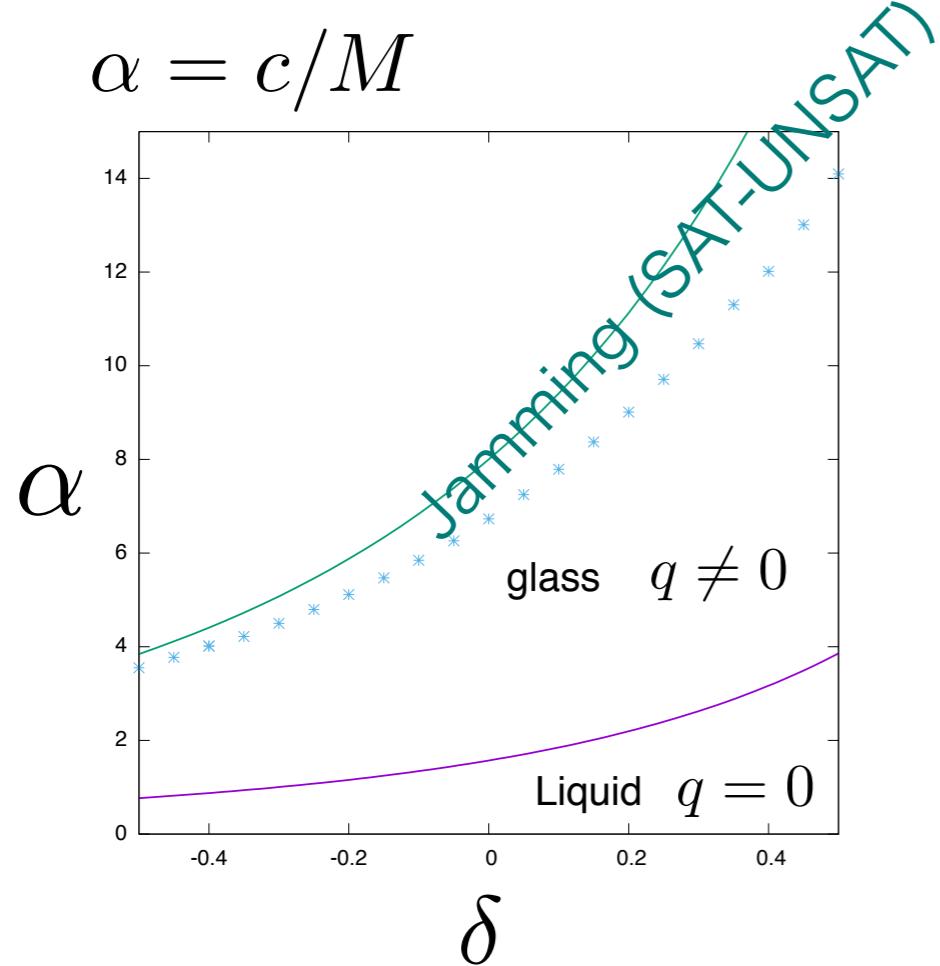
$$V(r) = \lim_{\epsilon \rightarrow \infty} \epsilon r^2 \theta(-r)$$

$M \rightarrow \infty$

H.Yoshino, SciPost Phys. 4 (6), 040 (2018)

continuous RSB:
hierarchical clustering of solutions

connectivity



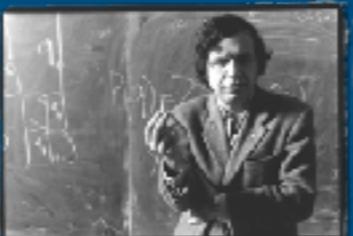
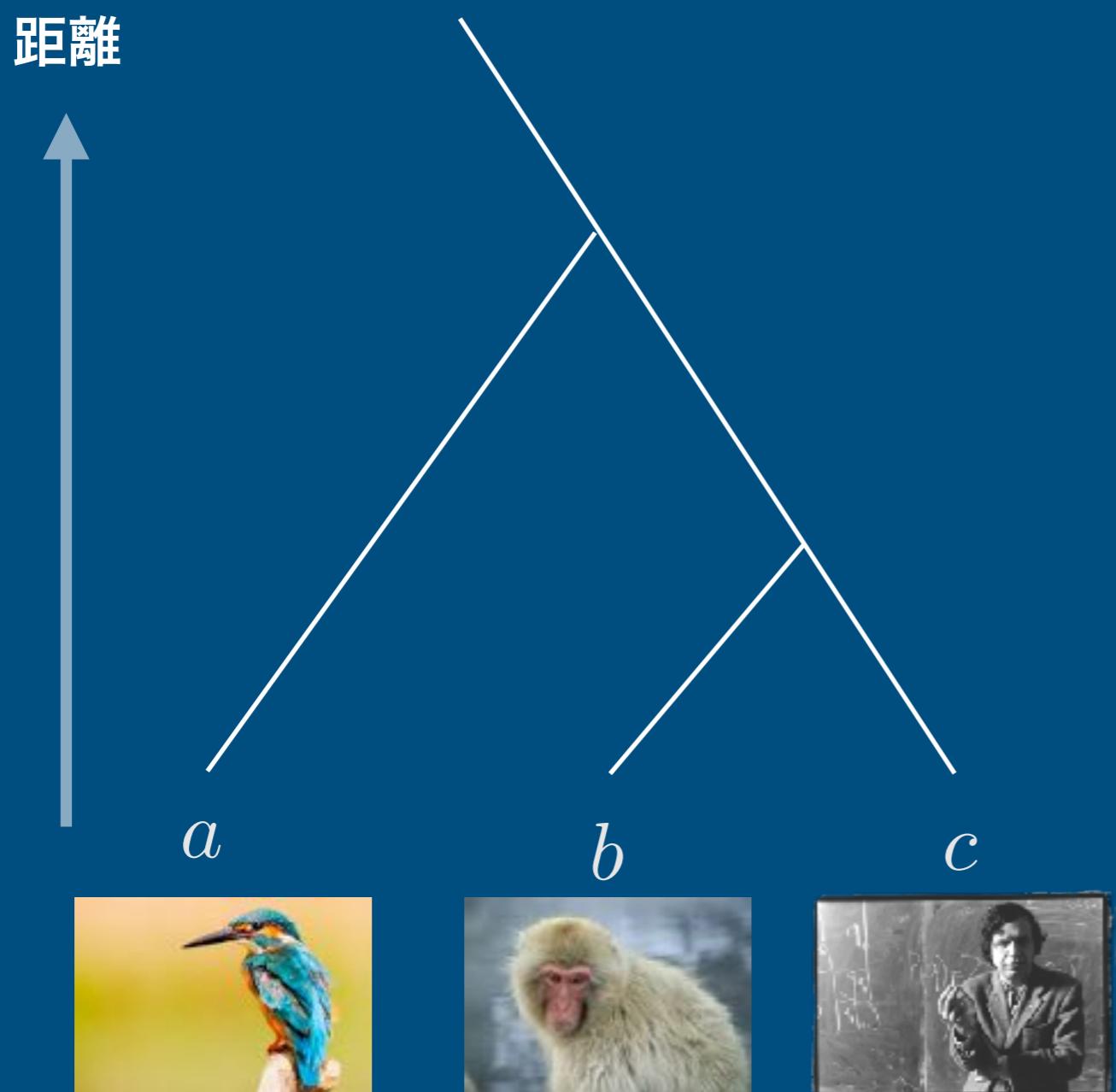
Same jamming criticality
as hard spheres
and perceptron ($p=1$)

Franz-Parisi (2016),
Franz-Parisi-Sevlev-
Urbani-Zamponi (2017)



Replica symmetry braking and ultra-metricity

a

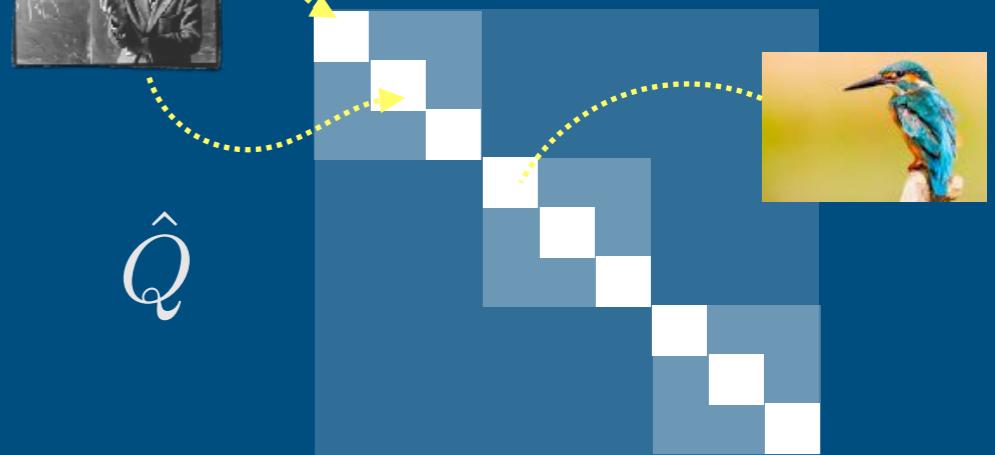


$$Q(a, b) = \min(Q(a, c), Q(b, c))$$



overlap matrix

\hat{Q}



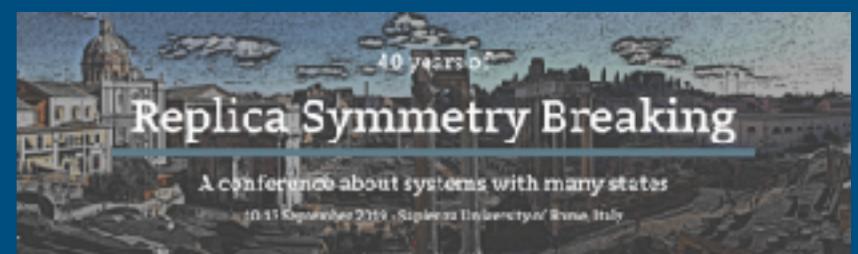
G. Parisi (1979)

first found in the
SK model for spin glass

Rigorous proof: M. Talagrand, (2003)

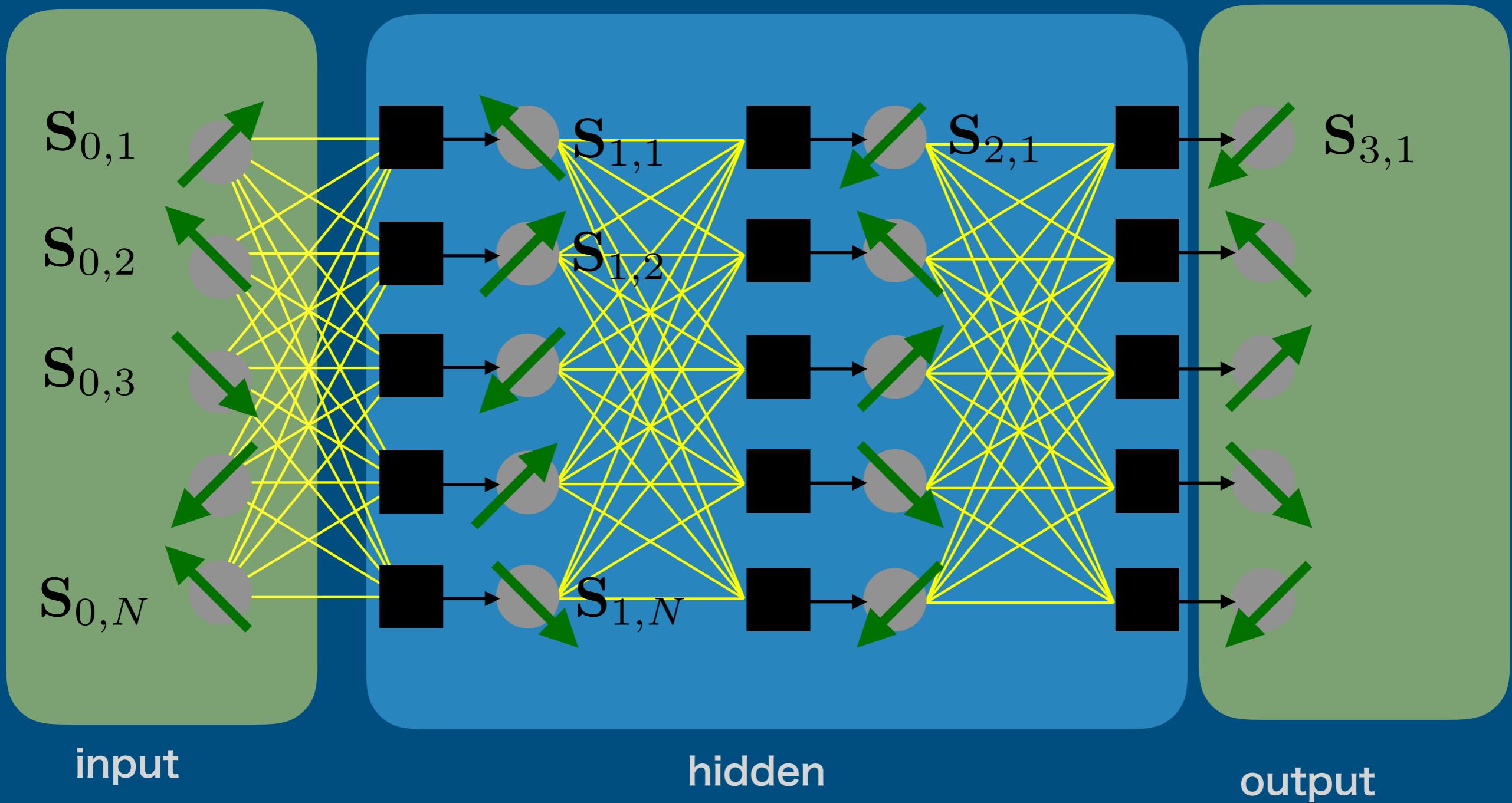
similarity

Q



Multi-layer Neural Network

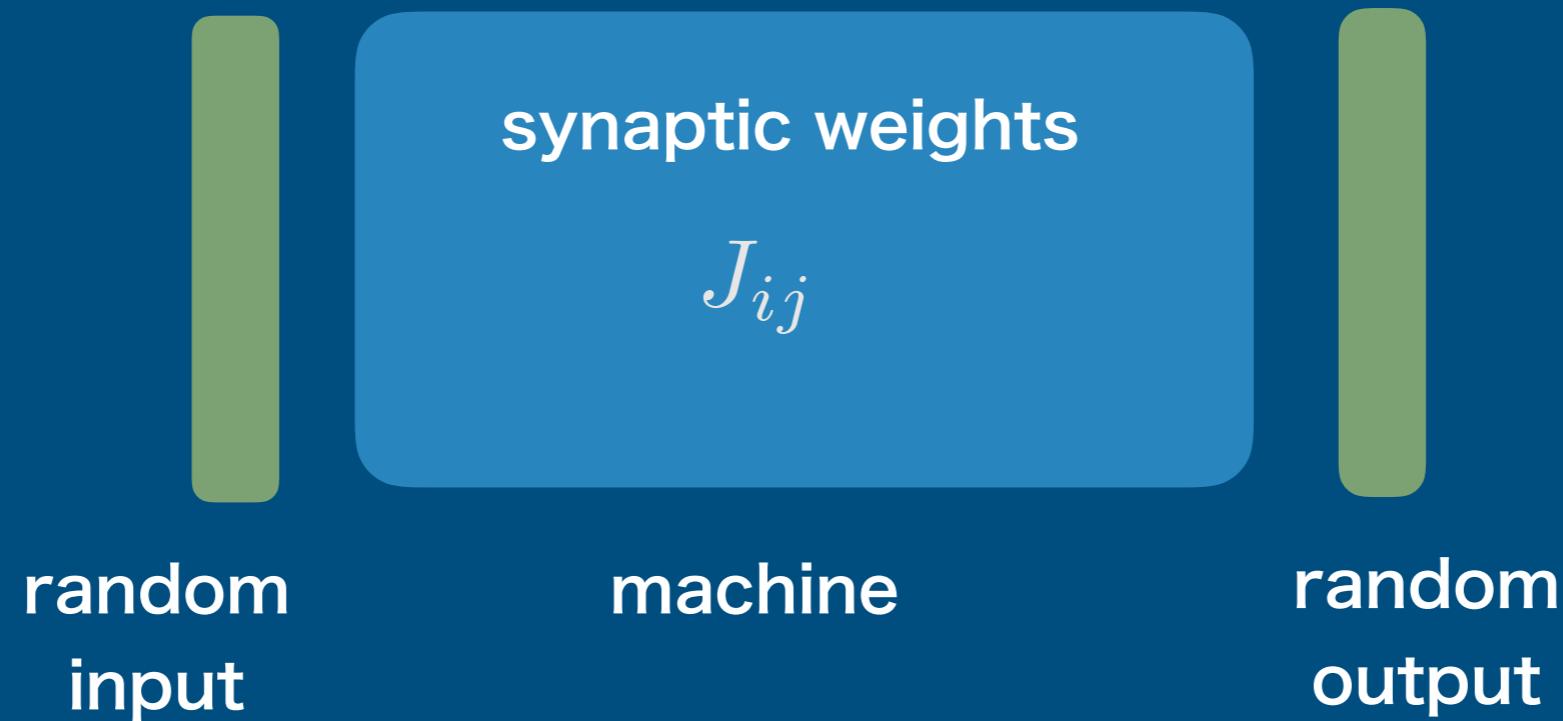
Design **weights** to satisfy **boundary conditions**



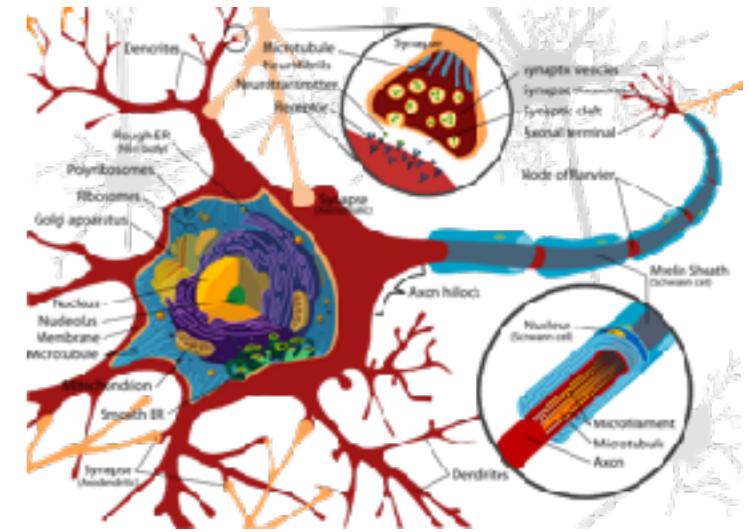
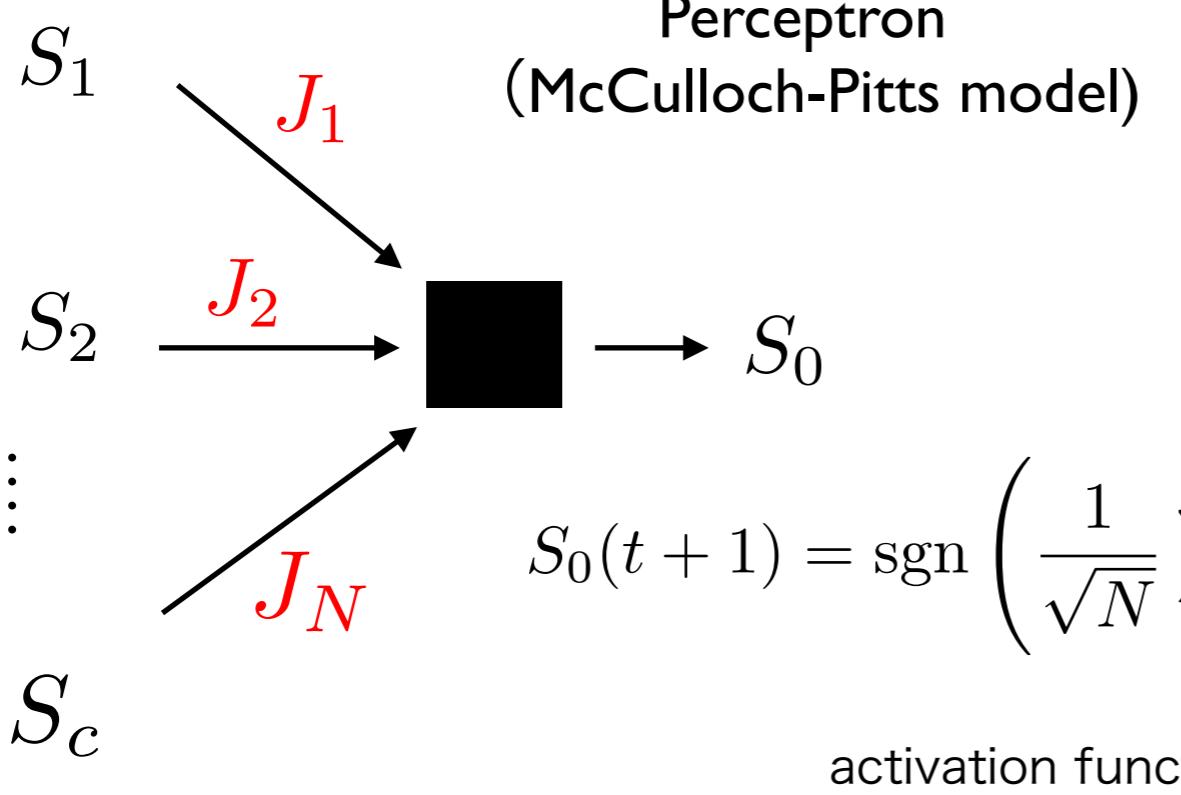
$\mathbf{S}_{i,l} = (S_{i,l}^1, S_{i,l}^2, \dots, S_{i,l}^M)$ $S_{i,l}^\mu = \pm 1$

Random inputs/random outputs

a constraint satisfaction problem (CSP)



Q: How many different ways the machine
can be designed to satisfy the
imposed random inputs/outputs ?



Elisabeth Gardner
(1957-1988)

Statistical mechanics on the “ensemble of fixed points”

fixed point patterns



$$S_i(t) = S_i^\mu \quad \mu = 1, 2, \dots, M$$

$$\alpha = \frac{M}{N} \quad M \rightarrow \infty \text{ with fixed } \alpha$$

Gardner volume

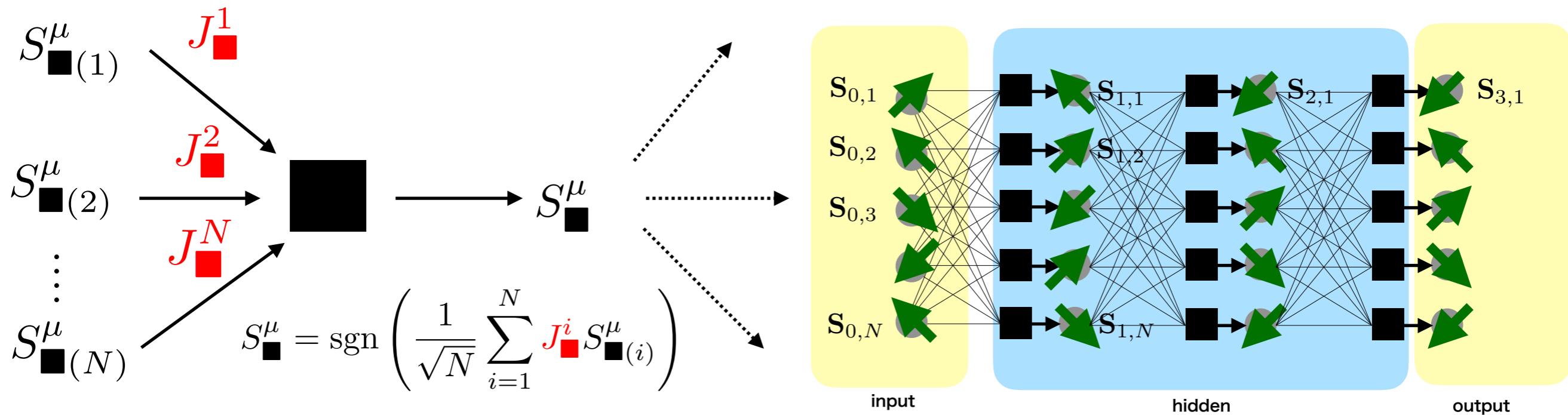
“Hardcore” constraint

$$e^{-\beta V(h)} = \theta(h)$$

$$V = \int \prod_{j=1}^N \frac{d\mathbf{J}_j}{\sqrt{2\pi}} e^{-\frac{\mathbf{J}_j^2}{2}} \prod_{\mu=1}^M e^{-\beta V(r^\mu)}$$

“Gap”

$$r^\mu = S_0^\mu \sum_{i=1}^N \frac{1}{\sqrt{N}} \mathbf{J}_i S_i^\mu$$



$$S_{L,i}^{\mu}(t+1) = \text{sgn} \left(\frac{1}{\sqrt{N}} \sum_{j=1}^N J_{L,i,j} \text{sgn} \left(\frac{1}{\sqrt{N}} \sum_{k=1}^N J_{L-1,j,k} \cdots \text{sgn} \left(\frac{1}{\sqrt{N}} \sum_{m=1}^N J_{1,l,m} S_{0,m}^{\mu}(t) \right) \right) \right)$$

Usual strategy of learning

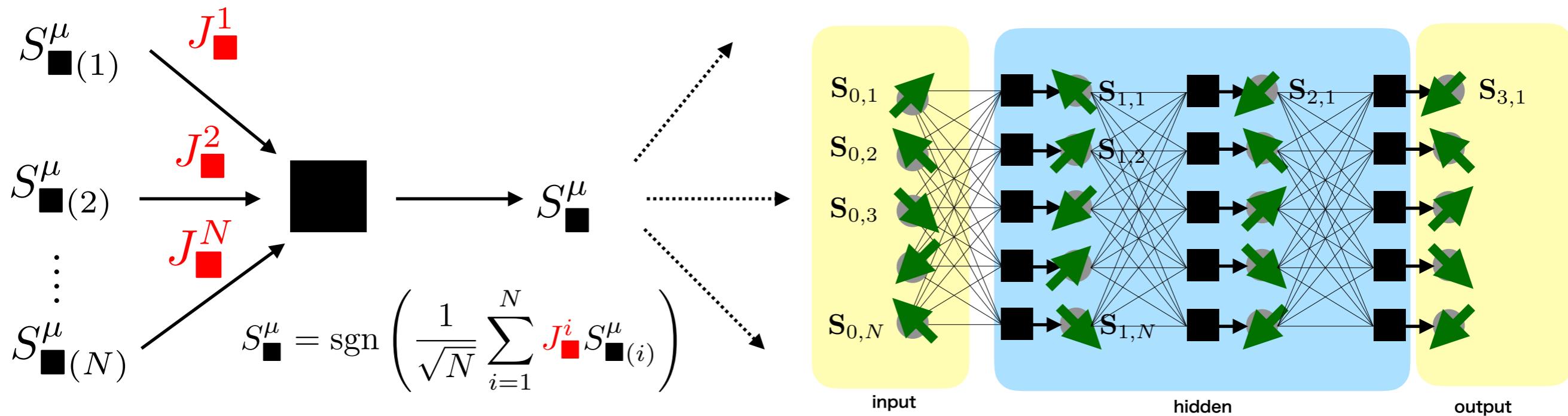
(1) define "loss function"

(2) try to minimize the loss function
via back-propagation

desired output
↓
e.g. $E = \sum_{i=1}^N \sum_{\mu=1}^M \left(S_{L,i}^{\mu} - (S_*)_{L,i}^{\mu} \right)^2$

e.g. SDG (stochastic gradient descent)

Too much long-ranged, highly convoluted, non-linear interaction! ...hard to analyze



Gardner volume generalized for a **multi-layer** network

trace over hidden variables

$$V(\mathbf{S}(0), \mathbf{S}(L)) = e^{NM\mathcal{S}(\mathbf{S}(0), \mathbf{S}(l))} = \left(\prod_{l=1}^{L-1} \prod_{i=1}^N \sum_{S_{l,i}^{\mu} = \pm 1} \right) \left(\int \prod_{\square} \prod_{j=1}^N \frac{dJ_{\square}^j}{\sqrt{2\pi}} e^{-\frac{(J_{\square}^j)^2}{2}} \right) e^{-\beta H}$$

**Hamiltonian with
“short-ranged” interactions**

$$H = \sum_{\mu=1}^M \sum_{\square} V(r_{\square}^{\mu})$$

“Hardcore” constraint

$$e^{-\beta V(h)} = \theta(h)$$

“Gap” $r_{\square}^{\mu} = \sum_{i=1}^N \frac{1}{\sqrt{N}} J^i S_{\square}^{\mu(i)} S_{\square}^{\mu}$

**Gaussian approx.
or modified model** $r_{\square}^{\mu} = \sum_{i=1}^N \frac{1}{\sqrt{N}} J^i \left(\frac{1}{\sqrt{M}} \sum_{\nu=1}^M \xi^{\mu\nu} S_{\square}^{\nu(i)} S_{\square}^{\nu} \right)$
 $\xi^{\mu\nu}$: Gaussian with zero mean and variance 1

(c.f.) internal representation (2-layer): R. Monasson and R. Zecchina (1995)

Replicated Gardner volume

replicas: machines learning in parallel

$$V^n(\mathbf{S}_0, \mathbf{S}_L) = \prod_{a=1}^n \left(\prod_{\blacksquare} \text{Tr}_{\mathbf{J}_{\blacksquare}^a} \right) \left(\prod_{\blacksquare \setminus \text{output}} \text{Tr}_{\mathbf{S}_{\blacksquare}^a} \right) \prod_{\mu, \blacksquare, a} e^{-\beta V(r_{\blacksquare, a}^\mu)}$$

$$r_{\blacksquare, a}^\mu = S_{\blacksquare, a}^\mu \sum_{i=1}^N \frac{1}{\sqrt{N}} J_{\blacksquare, a}^i S_{\blacksquare(i), a}^\mu$$

Glass order parameters

$$q_{ab, \blacksquare} = \frac{1}{M} \sum_{\mu=1}^M (S_{\blacksquare}^\mu)^a (S_{\blacksquare}^\mu)^b \quad Q_{ab, \blacksquare} = \frac{1}{N} \sum_{i=1}^N J_{\blacksquare(i)}^a J_{\blacksquare(i)}^b$$

Replicated free-energy

$$\frac{-\beta \overline{F(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} = \frac{\partial_n \overline{V^n(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} \Big|_{n=0} = S_n[\{\hat{Q}(l), \hat{q}(l)\}]$$

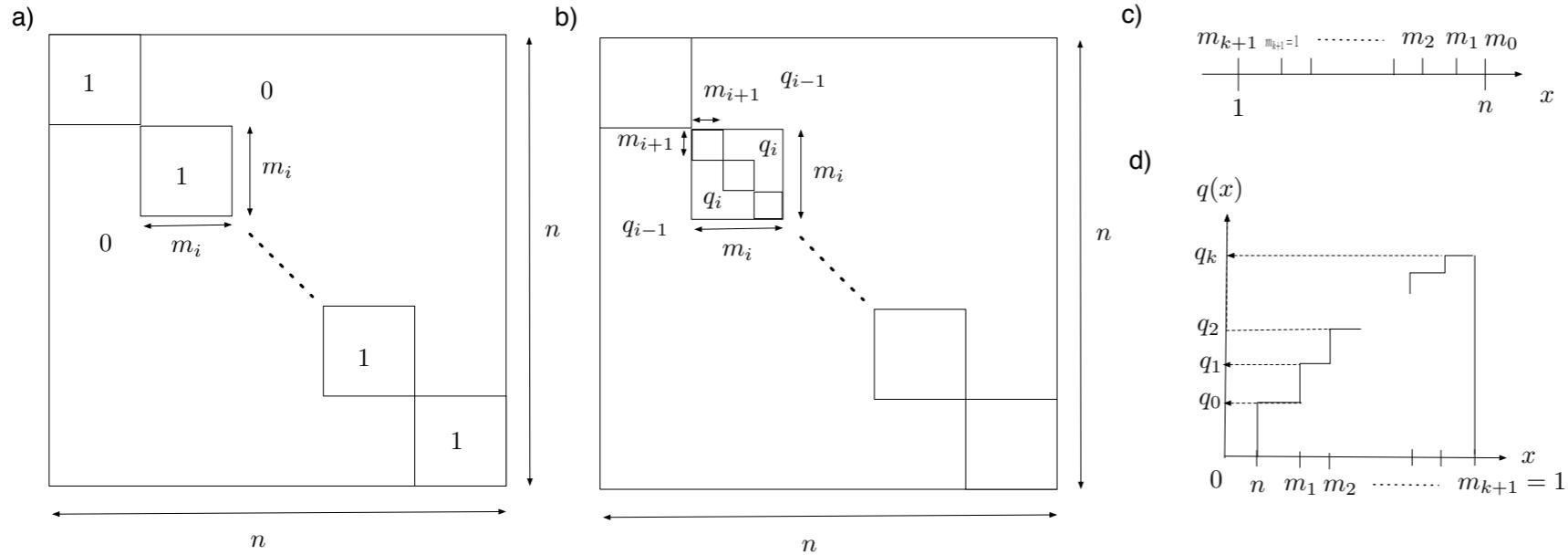
$$S_n[\{\hat{q}(l)\}, \{\hat{Q}(l)\}] = \alpha^{-1} \sum_{l=1}^L S_{\text{ent}}^{\text{bond}}[\hat{Q}(l)] + \sum_{l=1}^{L-1} S_{\text{ent}}^{\text{spin}}[\hat{q}(l)] \quad \alpha = \frac{M}{N}$$

$$- \sum_{l=1}^L e^{\frac{1}{2} \sum_{ab} q_{ab}(l-1) Q_{ab}(l) q_{ab}(l) \partial_{h_a(l)} \partial_{h_b(l)}} \prod_{a=1}^n e^{-\beta V(h_a(l))} \Big|_{h_a(l)=0}$$

quenched random input/output

$$q_{ab}(0) = q_{ab}(L) = 1$$

■ Parisi's RSB ansatz



overlap distribution function

$$P(q) = \frac{dx(q)}{dq}$$

$$\begin{aligned} Q_{ab}(l) &= \sum_{i=0}^{k+1} Q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L \\ q_{ab}(l) &= \sum_{i=0}^{k+1} q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L-1 \\ \varepsilon_{ab}(l) &= \sum_{i=0}^k \varepsilon_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L-1 \end{aligned}$$

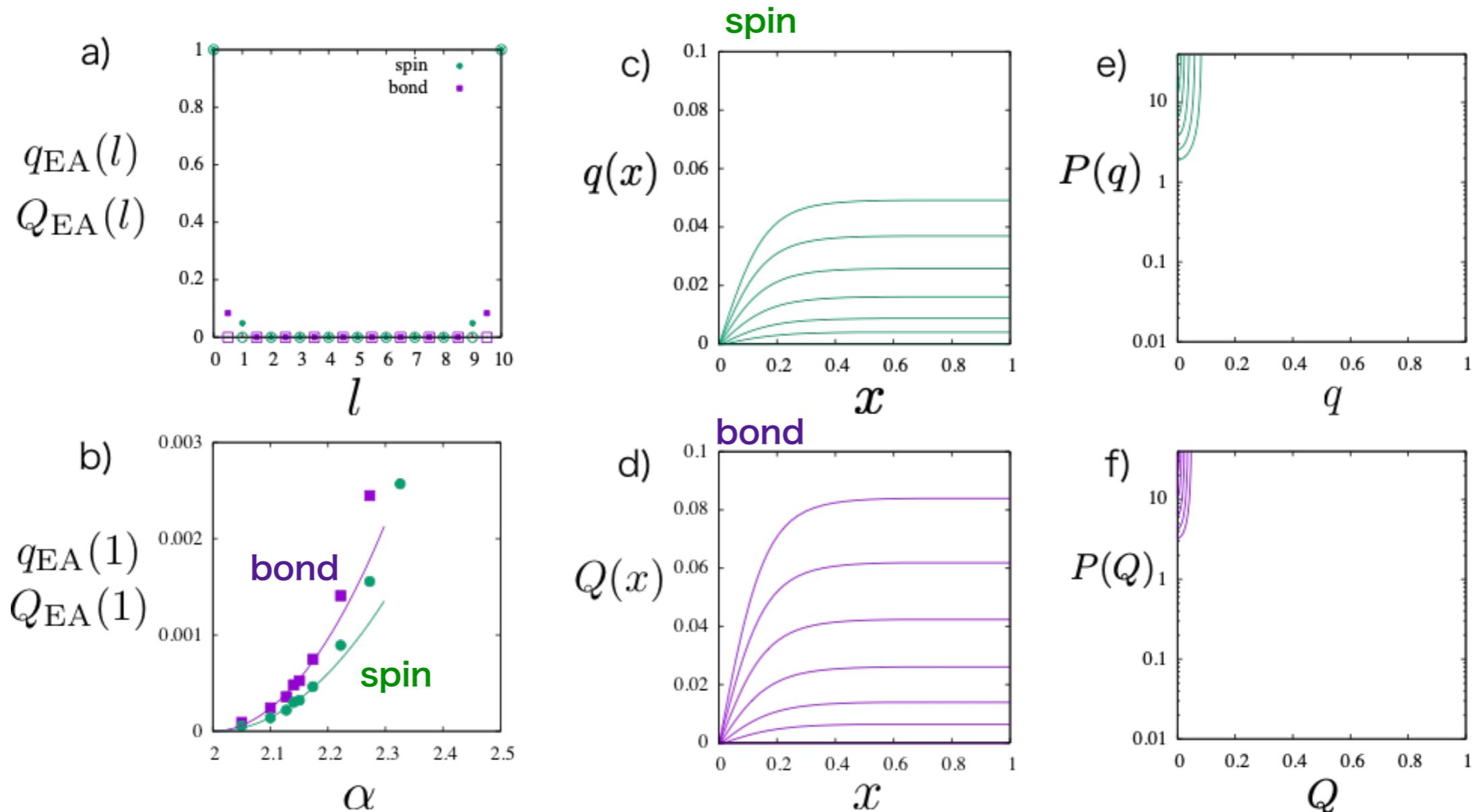
entropic part (bond)

$$S_{\text{ent}}[\hat{Q}(l)] = \frac{1}{2} \ln \det \hat{Q}(l)$$

entropic part (spin)

$$S_{\text{ent}}[\hat{q}(l)] = \sum_{a < b} \ln e^{-\sum_{a < b} \varepsilon_{ab}^*(l) \frac{\partial^2}{\partial h_a \partial h_b}} \prod_a 2 \cosh(h_a) \Bigg|_{\{h_a=0\}} \quad q_{ab} = -\frac{\delta}{\delta \varepsilon_{ab}} \ln e^{-\sum_{a < b} \varepsilon_{ab} \frac{\partial^2}{\partial h_a \partial h_b}} \prod_a 2 \cosh(h_a) \Bigg|_{\{h_a=0\}} \Bigg|_{\varepsilon_{ab}=\varepsilon_{ab}^*[\hat{q}]}$$

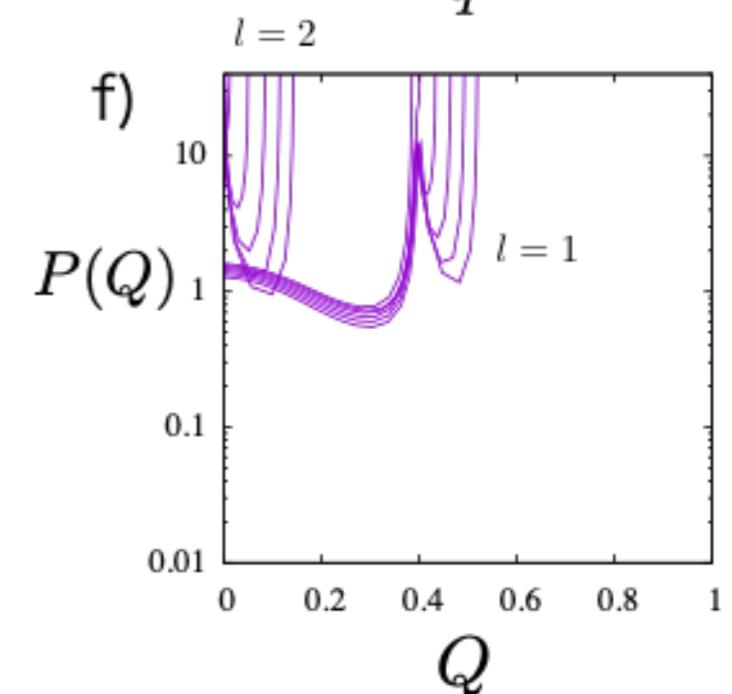
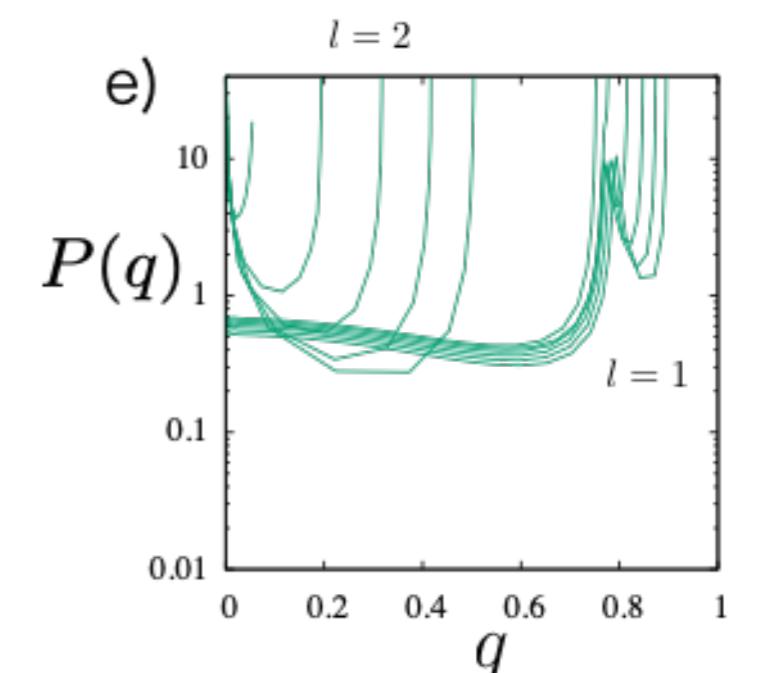
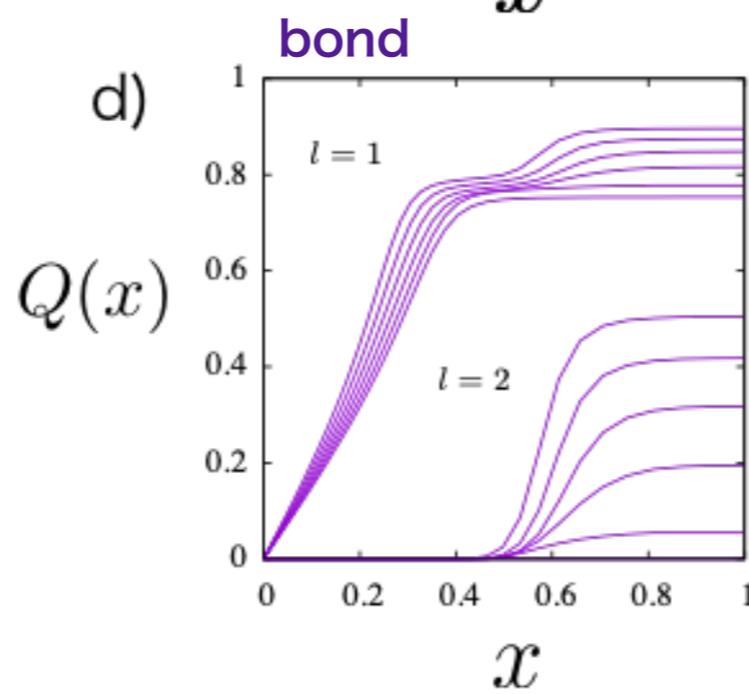
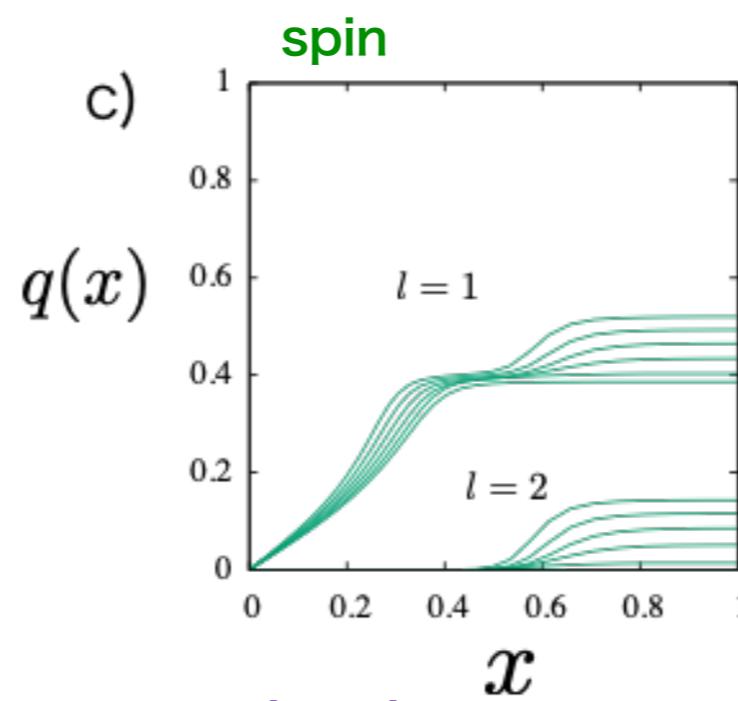
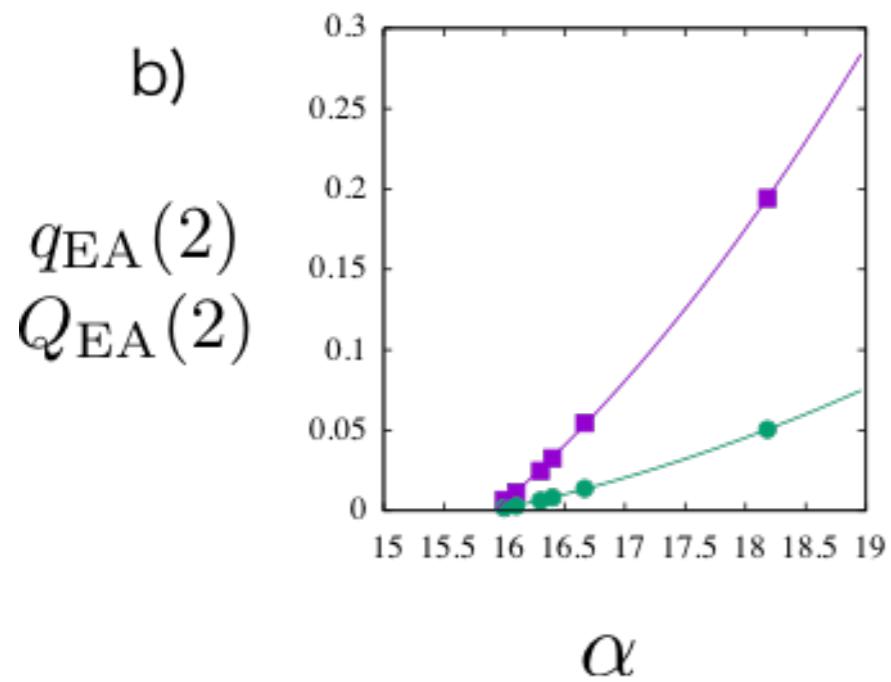
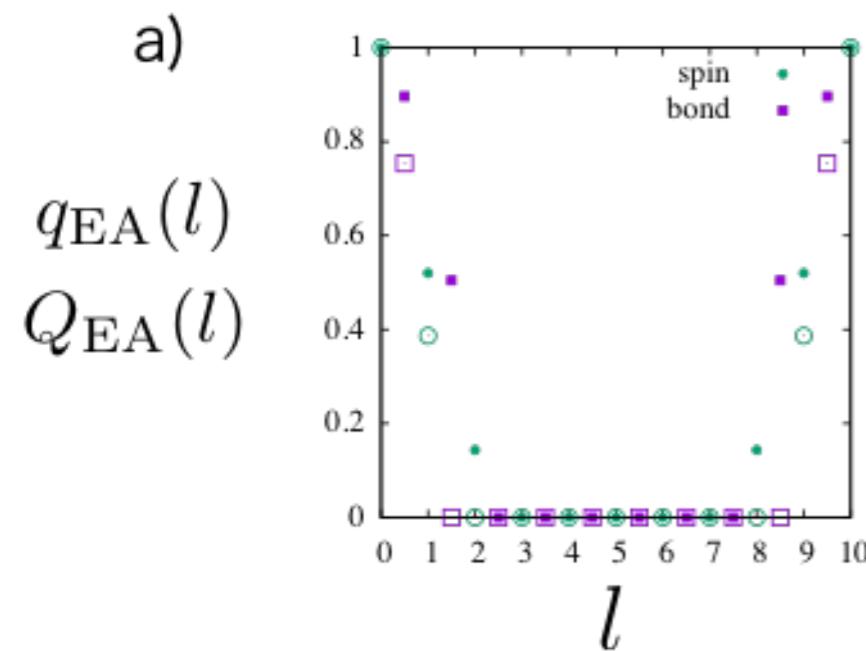
■ 1st Glass transition



$$\alpha_g \simeq 2.03$$

continuous transition to full RSB glass phase
at 1 st & ($L-1$) th layer
other layers remain in the liquid phase

■ 2nd Glass transition

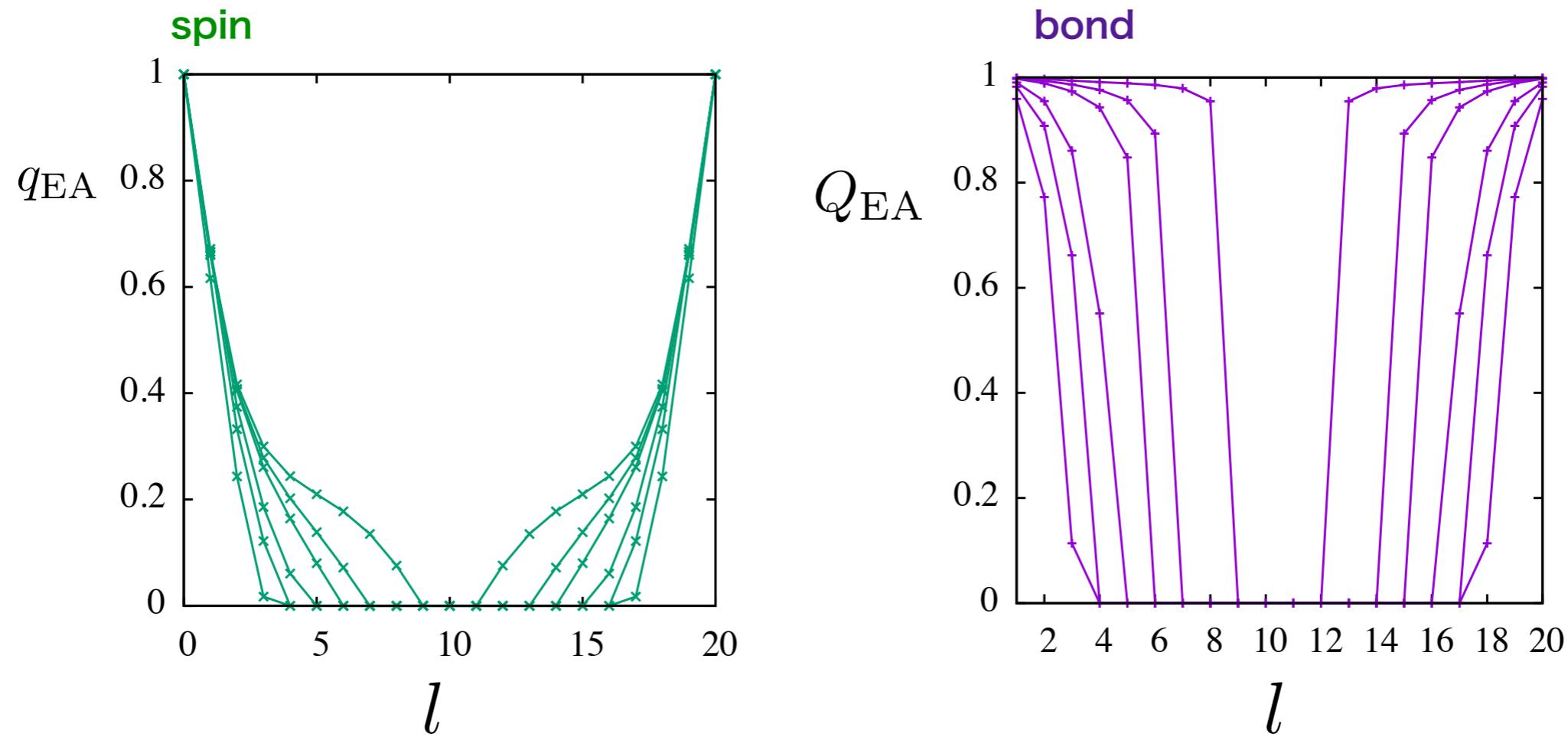


$$\alpha_g(2) \simeq 15.38$$

continuous transition to full RSB glass phase
at 2nd & (L-2) th layer

which also induce 2nd glass transitions at 1st and L-th layer

■Growth of glass phase under larger constrains



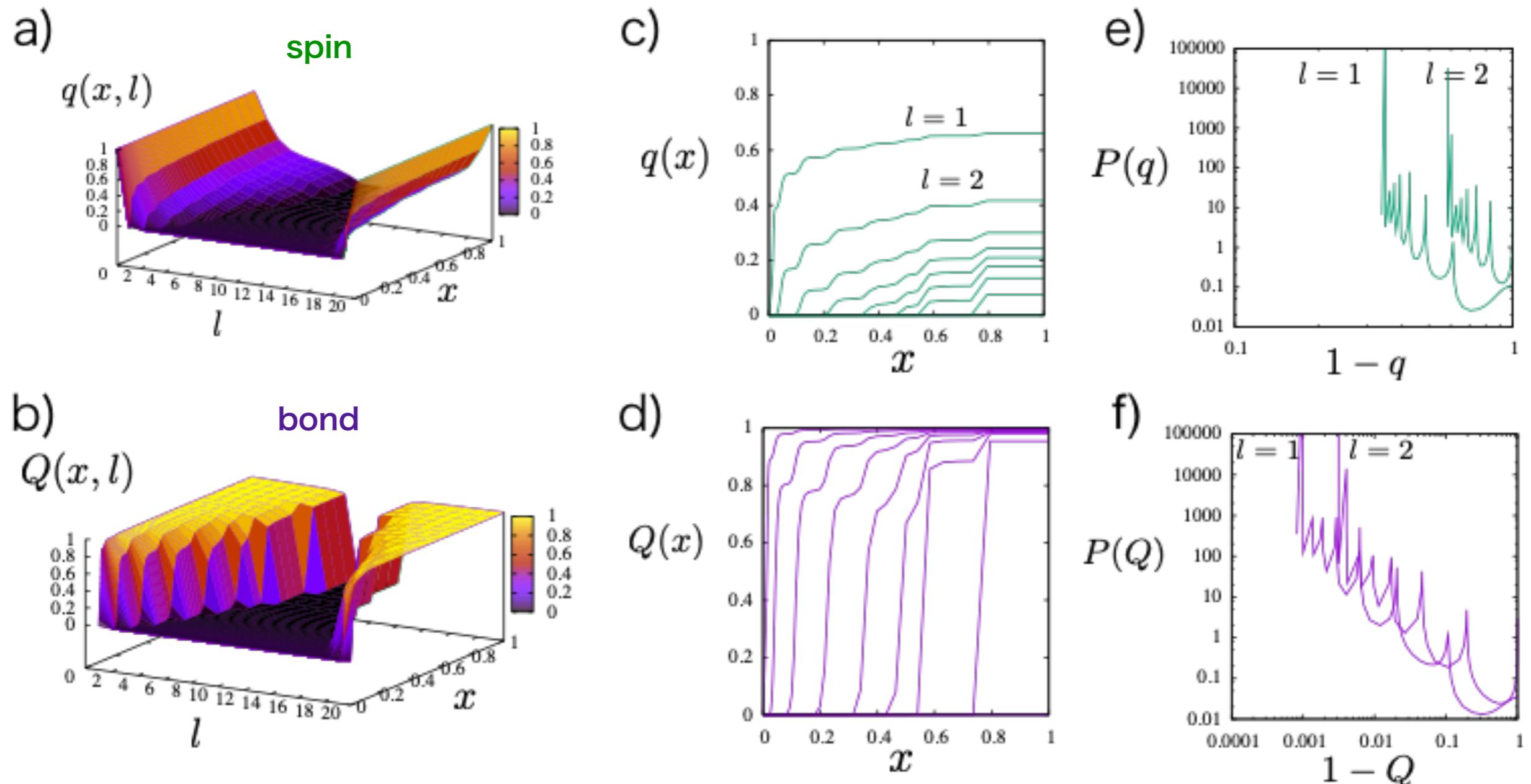
$$1/\alpha = 0.02, 0.01, 0.005, 0.001, 0.0005, 0.00025$$

“penetration depth”

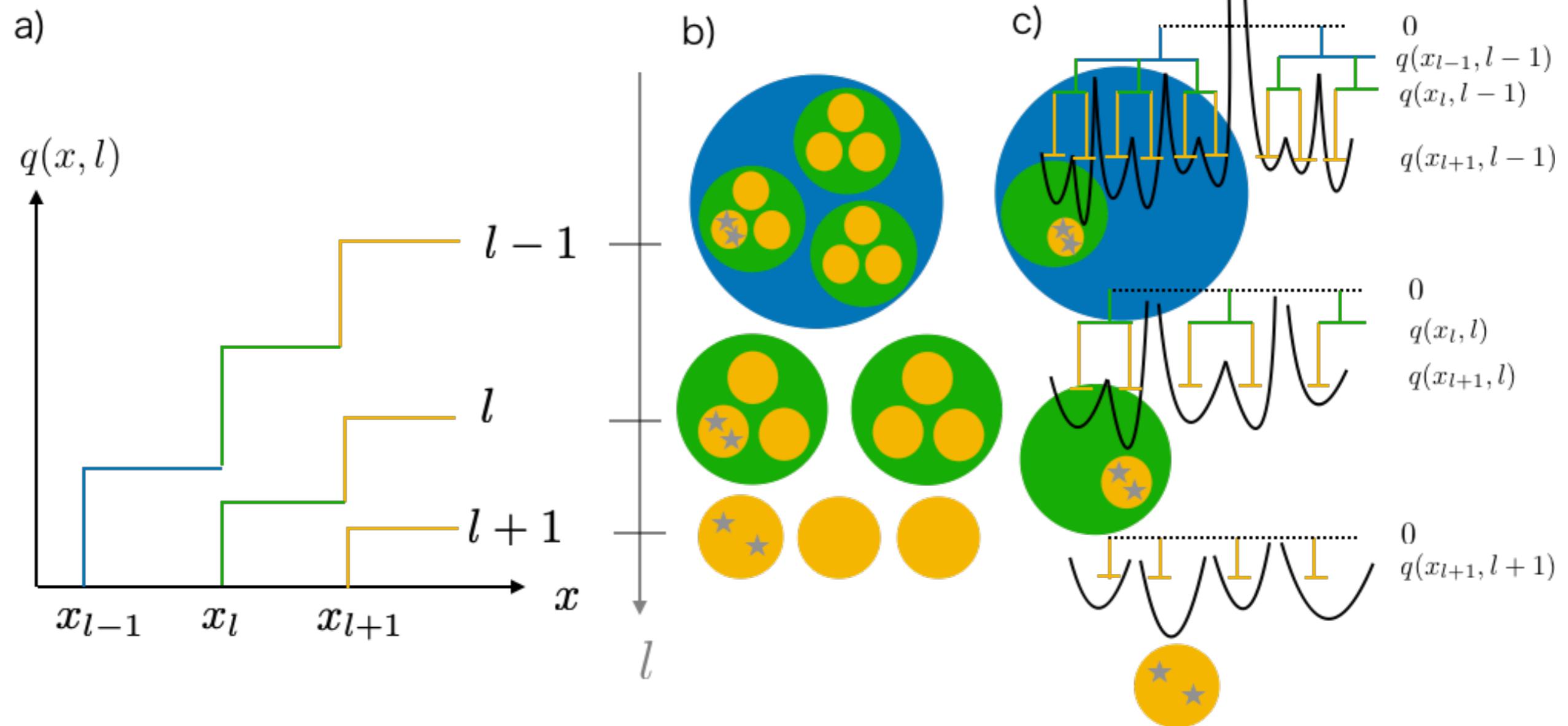
$$\xi(\alpha) \propto \ln \alpha$$

■ More “terraces” under larger constraints

$$\alpha = 4000$$



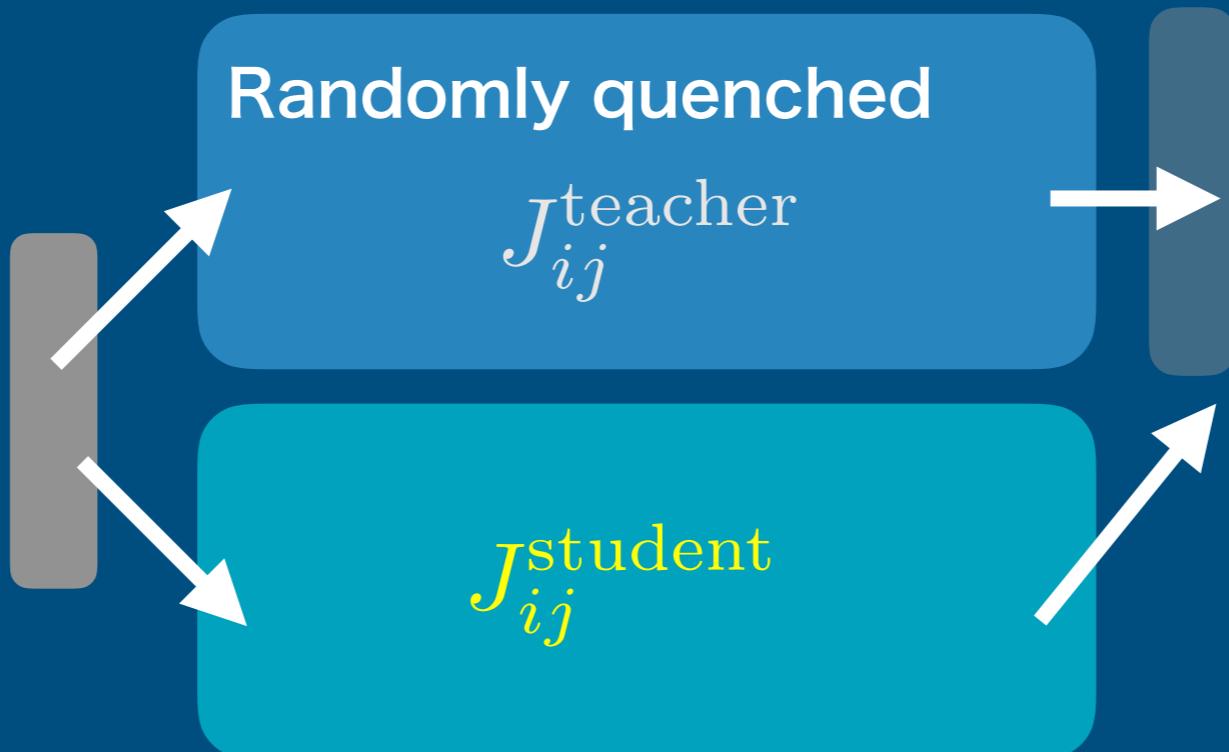
■ Summary: depth dependent free-energy landscape



Teacher student setting a statistical inference problem (SIP)

1) Training

random
test
data

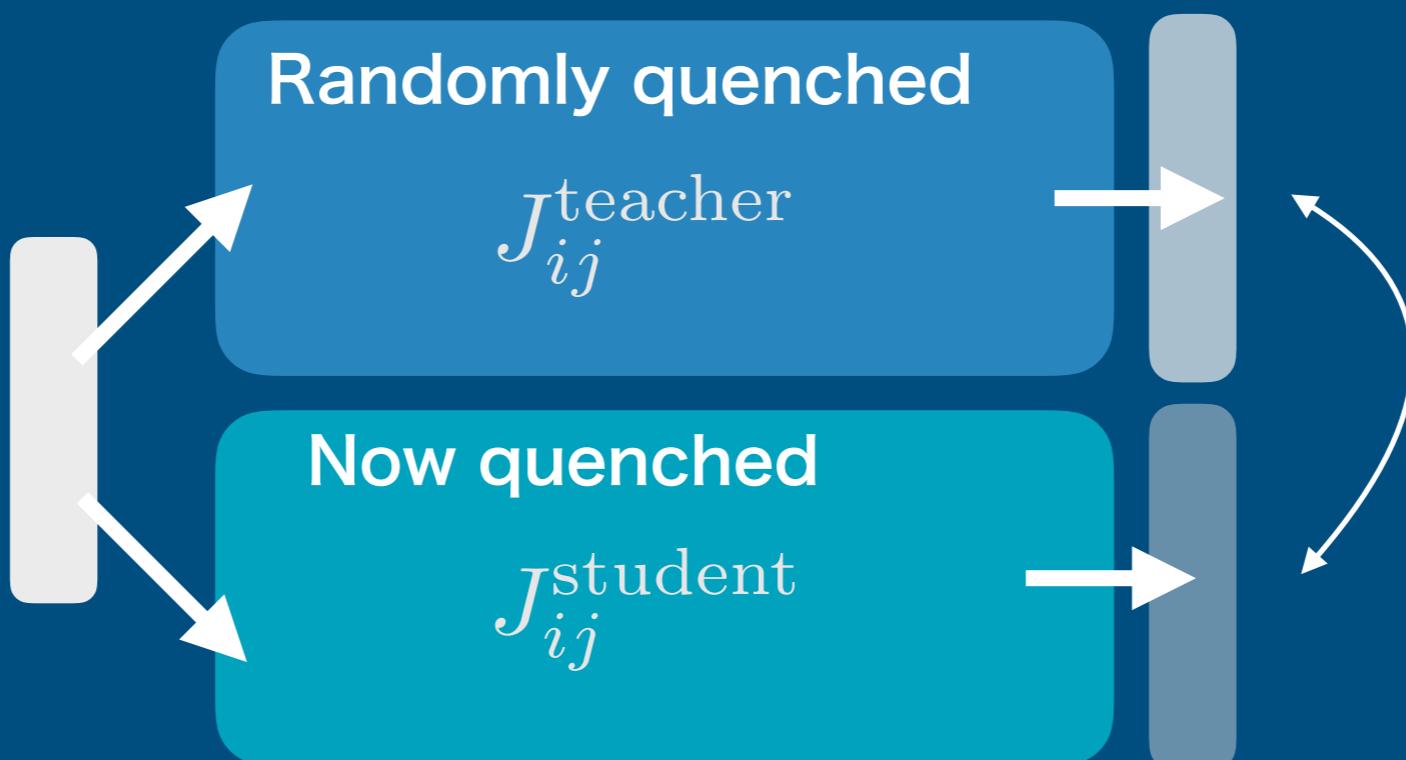


Out put of teacher

Student is forced to
reproduce
teacher's output

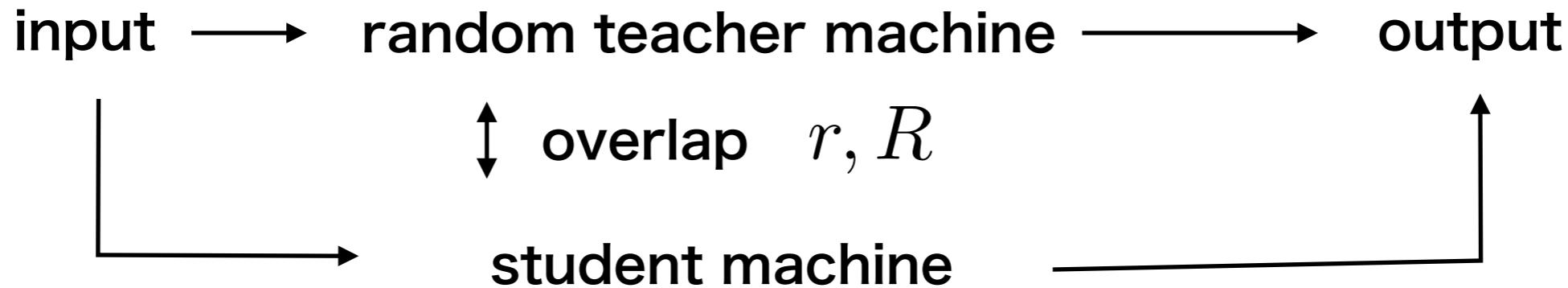
2) Test

random
test
data



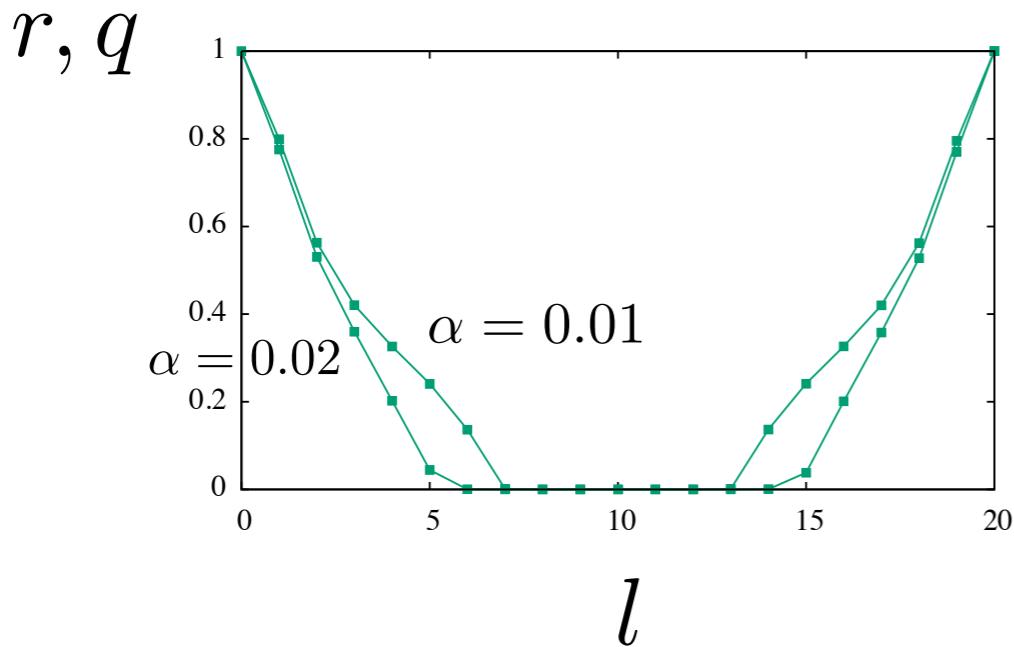
Compare

teacher+student machine (inference)



Review: Zdeborova-Krzakala (2017)

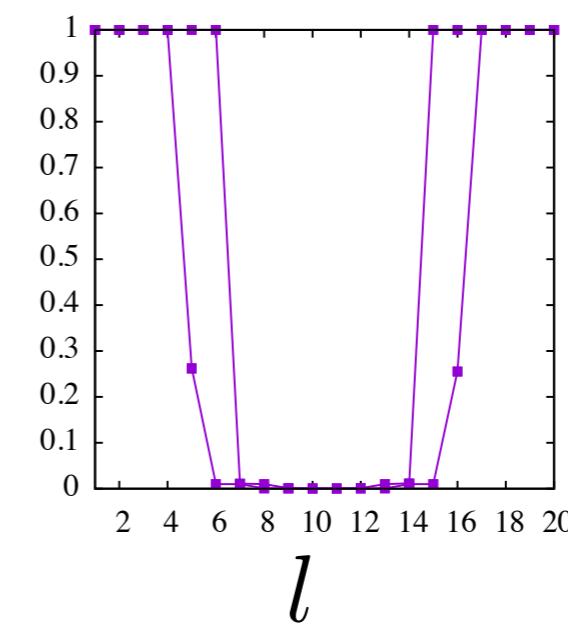
I+s replica



RS solution

Bays optimal, Nishimori condition

R, Q

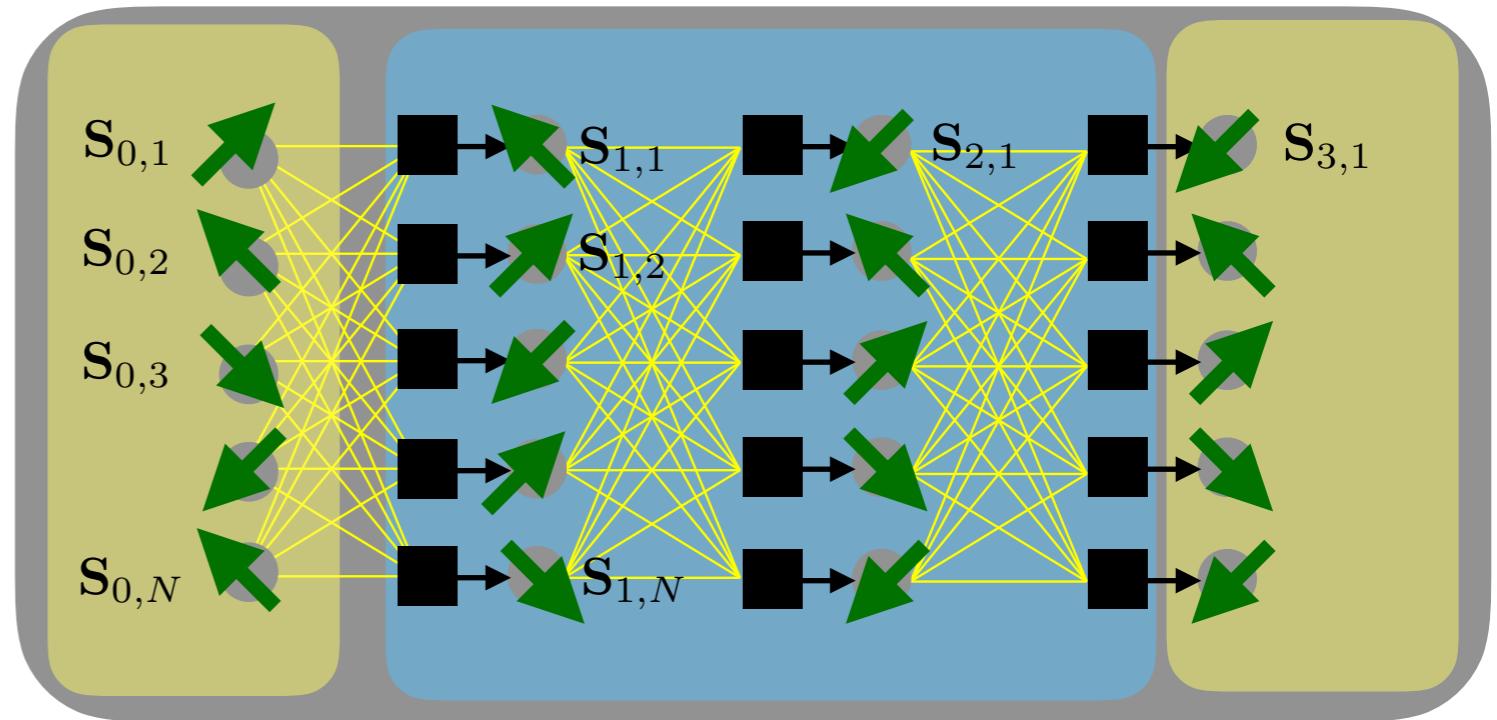


Symmetry breaking field: remanent bias in the liquid phase

$O(\ln(N)/N)$

Simulations of learning

■ random inputs/outputs



Monte Carlo simulation

Hamiltonian

$$H = - \sum_{\blacksquare} \sum_{\mu=1}^M V(r_{\blacksquare}^\mu)$$

Gap

$$r_{\blacksquare}^\mu = \sum_{i=1}^N \frac{1}{\sqrt{N}} J_{\blacksquare}^i S_{\blacksquare(i)}^\mu S_{\blacksquare}^\mu$$

soft-core potential

$$e^{-\beta V(h)} = e^{-\beta \epsilon h^2} \theta(h) \quad \beta = \frac{1}{k_B T}$$

Dynamical variables

$$J_{\blacksquare}^i, S_{\blacksquare}^\mu \quad (i = 1, 2, \dots, N) (\mu = 1, 2, \dots, M) \\ (\blacksquare = 1, \dots, LN)$$

with random boundaries

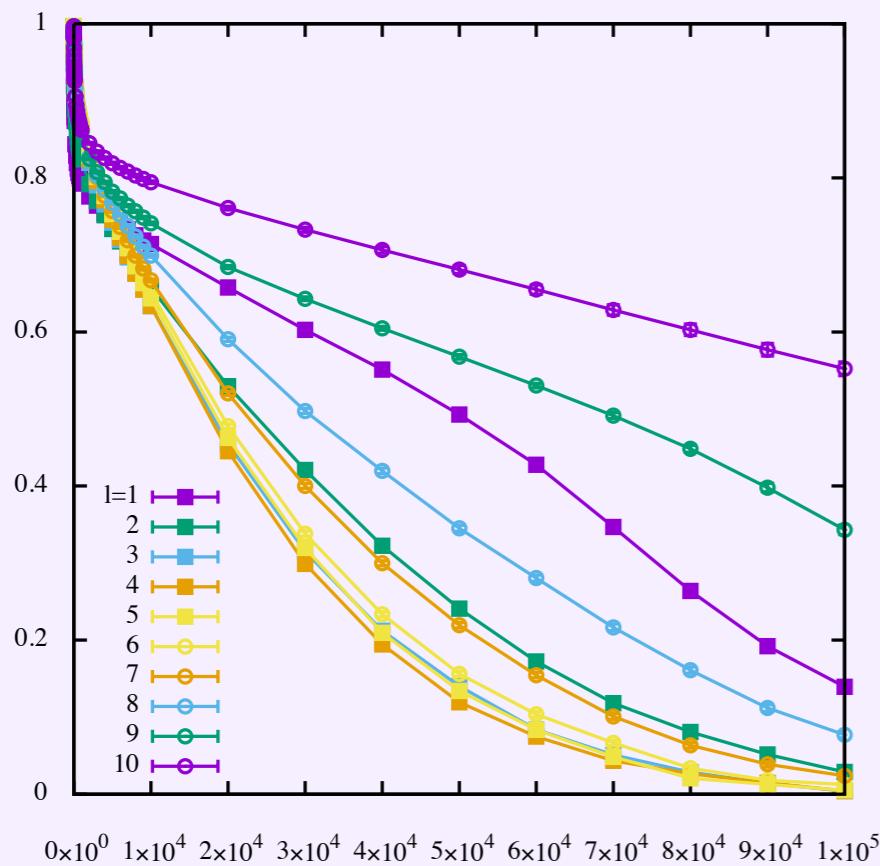
$$S_{0,i}, S_{L,i} \quad (i = 1, 2, \dots, N)$$

Relaxation of autocorrelation functions

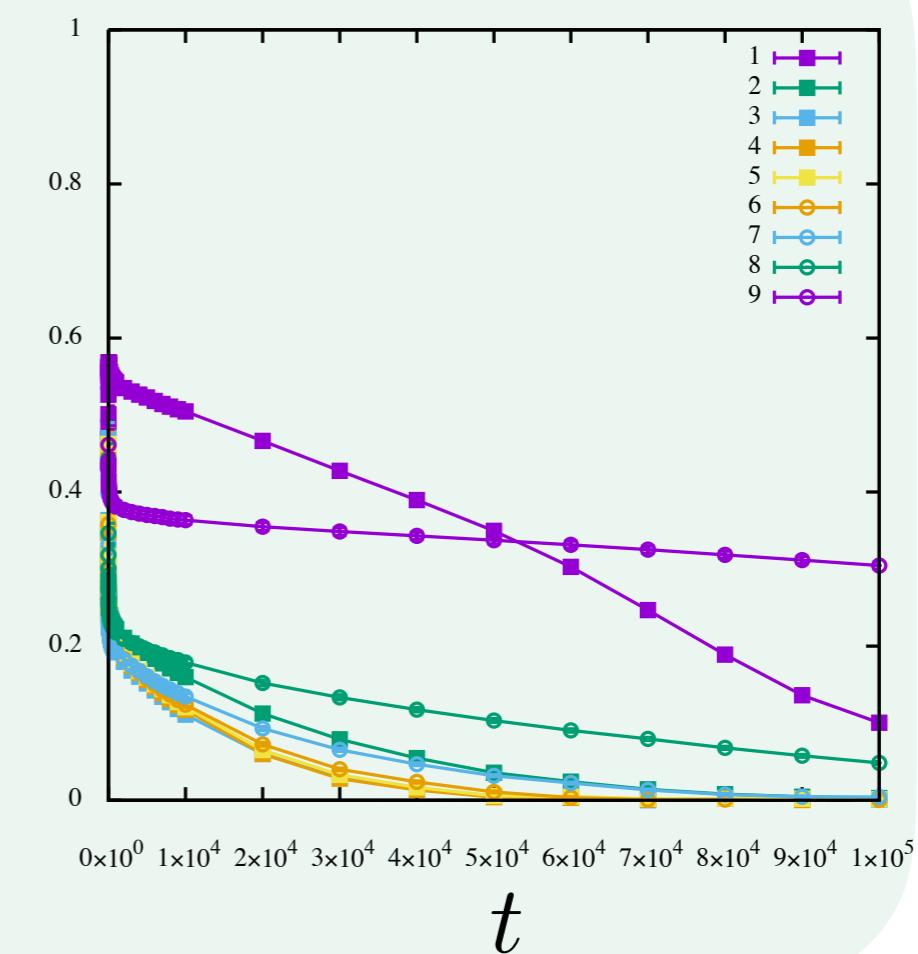
$$C_{\text{bond}}(t, \blacksquare) = \frac{1}{N} \sum_{i=1}^N \langle J_{\blacksquare}^i(0) J_{\blacksquare}^i(t) \rangle$$

$$C_{\text{spin}}(t, \blacksquare) = \frac{1}{M} \sum_{\mu=1}^M \langle S_{\blacksquare}^{\mu}(0) S_{\blacksquare}^{\mu}(t) \rangle$$

$C_{\text{bond}}(t, l)$



$C_{\text{spin}}(t, l)$



$N = 20, M = 200 (\alpha = 10)$

$L = 10 \quad T = 0.015 \quad 240 \text{ samples}$

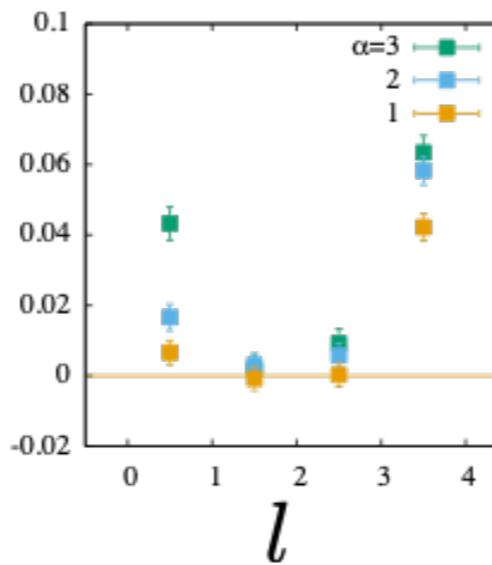
■ teacher-student setting

binary perceptron

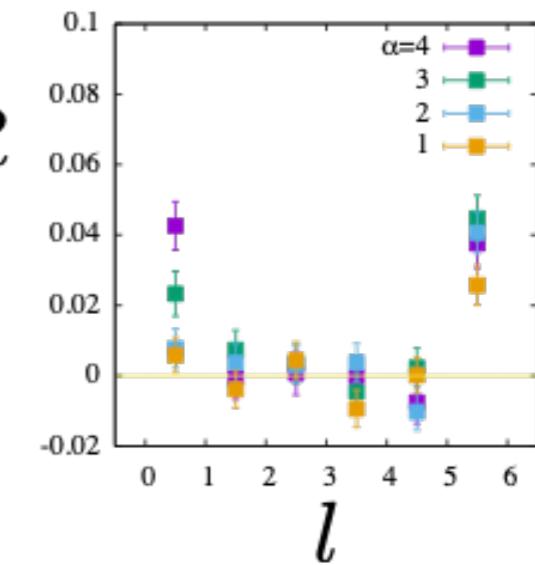
training

a)
bond

R



R



test

b)
spin

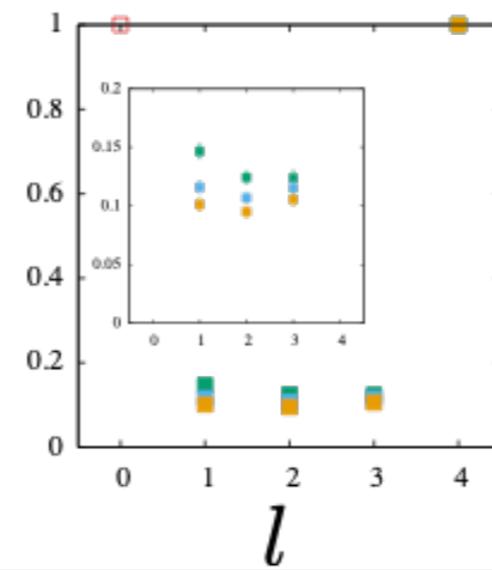
r_{test}

c)

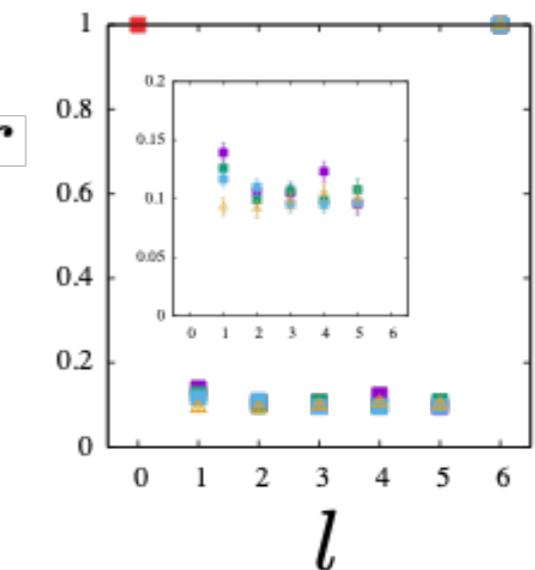
spin

Overlaps between
teacher and student

r



r

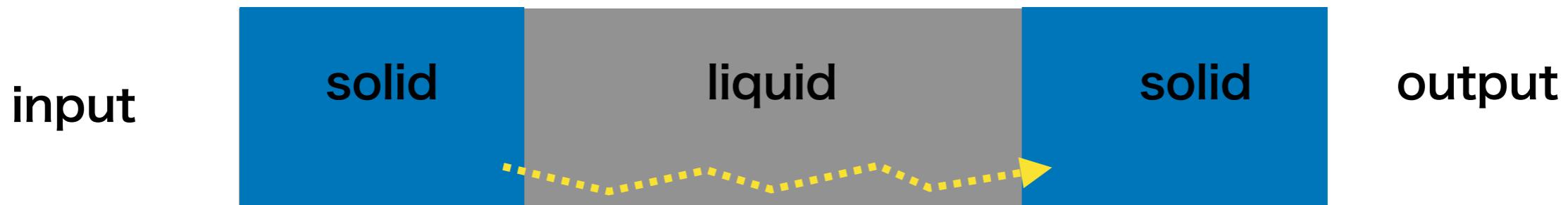


r_{test}

l

Summary

The solution space of over-parametrized DNN exhibit “solid-liquid-solid” structure



1. Learning : “Liquid phase” helps equilibration
2. Generalization: “crystal phase” enables generalization. Remanent bias field in the liquid phase plays the role of symmetry breaking field by which “hidden” crystal is selected out of many glass solutions.
3. Space evolution of the hierachal free-energy landscape: DNN naturally has the power of renormalization: classification, feature detection

Outlook

Numerical simulations to test theoretical predictions

Various statistical inference problems



Complex systems with heterogeneity

- gene regulatory network, ...
functionality vs robustness in biology
- allostericity
- ultra-stable glass, rheology

Theories in $d = 1 + \infty$

Like “Landau to Ginzburg-Landau” but more microscopic