

IPISeminar Oct 28, 2021

Design space of a deep neural network - its spatial evolution and robustness

Hajime Yoshino

Cybermedia Center, Osaka University



Hajime Yoshino, SciPostPhys. Core 2, 005 (2020).

「最近の研究から - 深層ニューラルネットワークの解剖—統計力学によるアプローチ」日本物理学会誌76巻9号(2021年9月号)

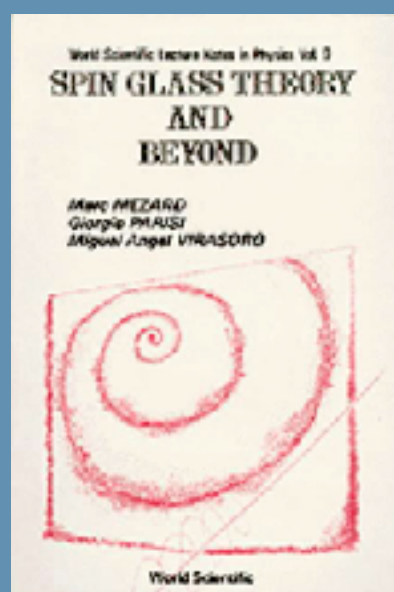


1978

“Deep Learning”

statistical mechanics of disordered systems

without quenched disorder



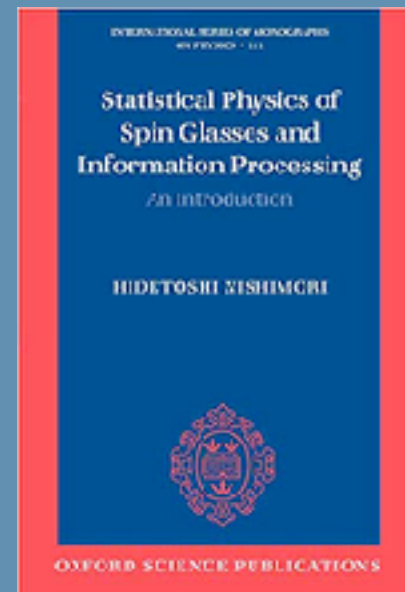
1987



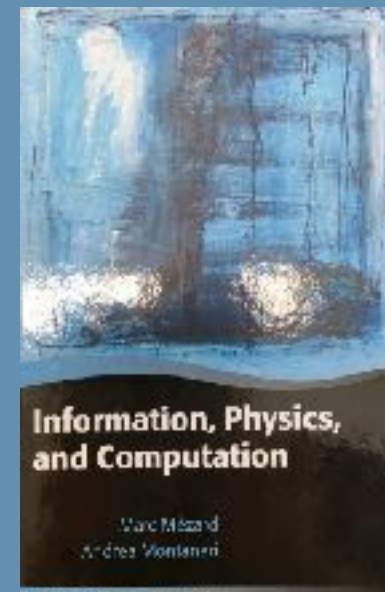
1991



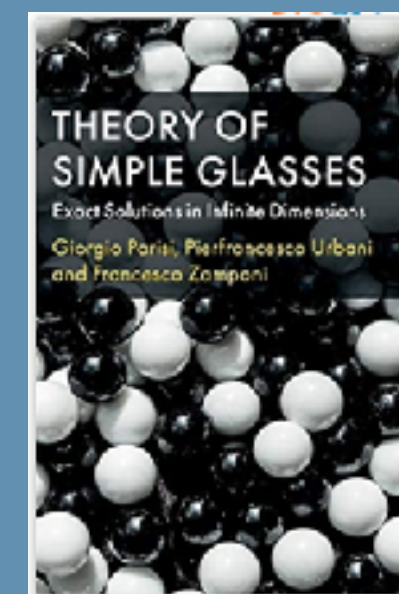
1999



2001



2008



2020

spin glass theory and FAR beyond

in progress

Statistical mechanics of disordered systems: spins, spheres and machines

Hajime Yoshino^{1,2}

¹Cyberrena Center, Osaka University, Toyonaka, Osaka 565-0843, Japan

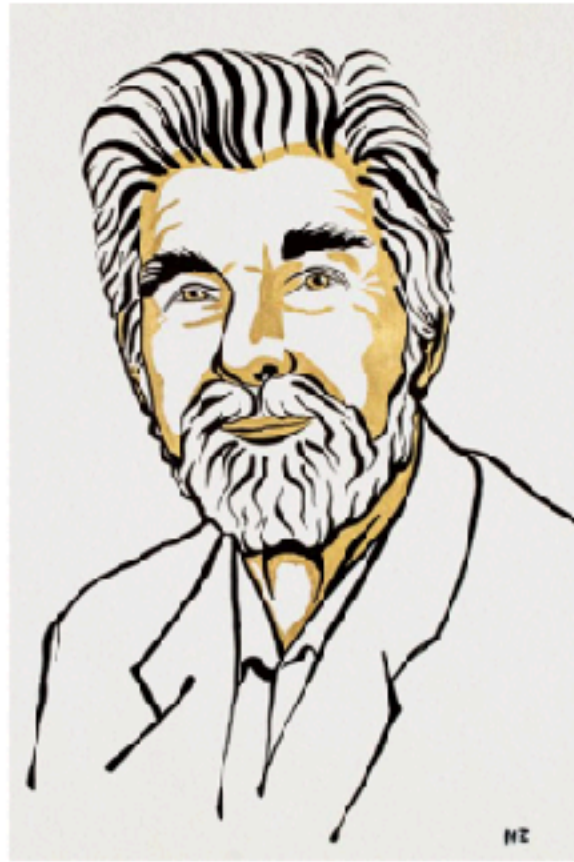
²Graduate School of Science, Osaka University, Toyonaka, Osaka 565-0843, Japan

In this lecture note we discuss glass physics and related problems using solvable mean-field models. First we discuss mean-field spin models without quenched disorder (just ferromagnetic couplings) with dense (but not global) couplings. We show that they exhibit glassy phases in supercooled paramagnetic phase and recover the standard results known in the mean-field spin-glass models with quenched disorder. Next we discuss glass physics in dense assemblies of simple spheres in large dimensional limit.

The Nobel Prize in Physics 2021



III. Niklas Elmehed © Nobel Prize Outreach
Syukuro Manabe
Prize share: 1/4



III. Niklas Elmehed © Nobel Prize Outreach
Klaus Hasselmann
Prize share: 1/4



III. Niklas Elmehed © Nobel Prize Outreach
Giorgio Parisi
Prize share: 1/2



https://www.nobelprize.org/uploads/2021/10/sciback_fy_en_21.pdf

“for groundbreaking contributions to our understanding of complex physical systems”



Perceptron

Example



of possible machines
 $2^3 = 8$

machine1

machine2

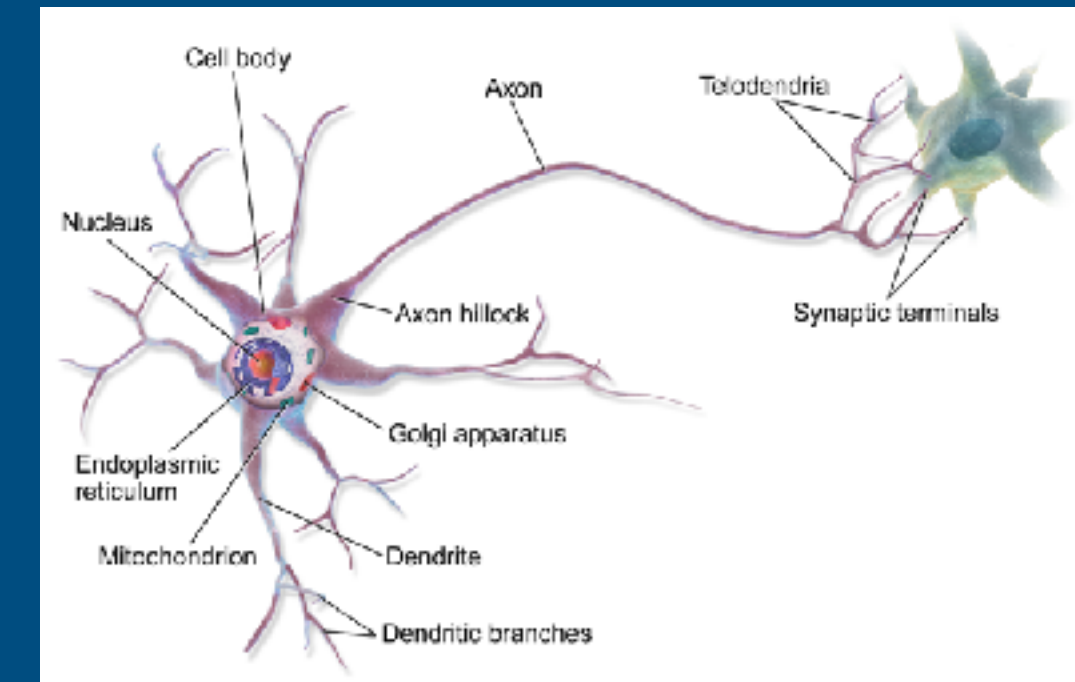
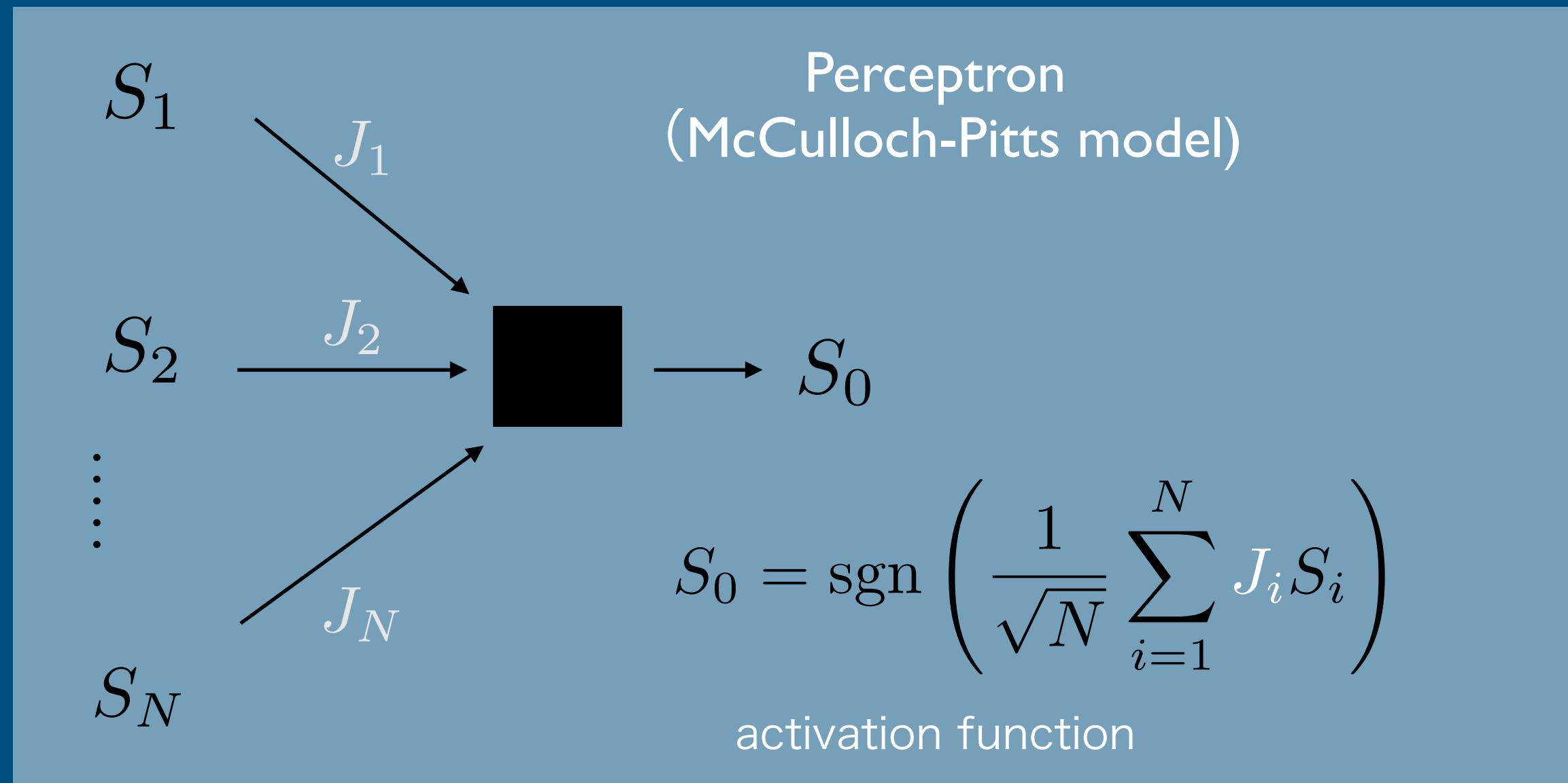
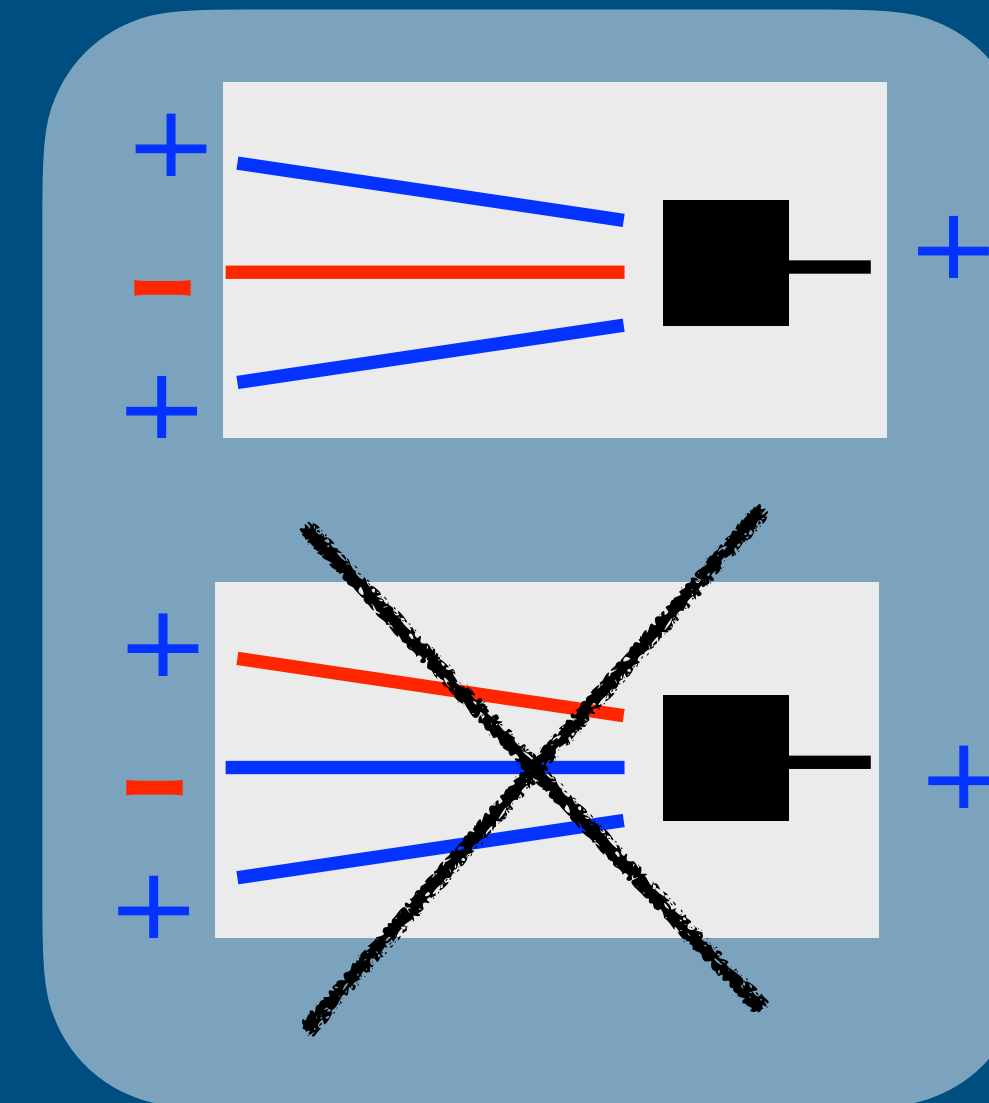
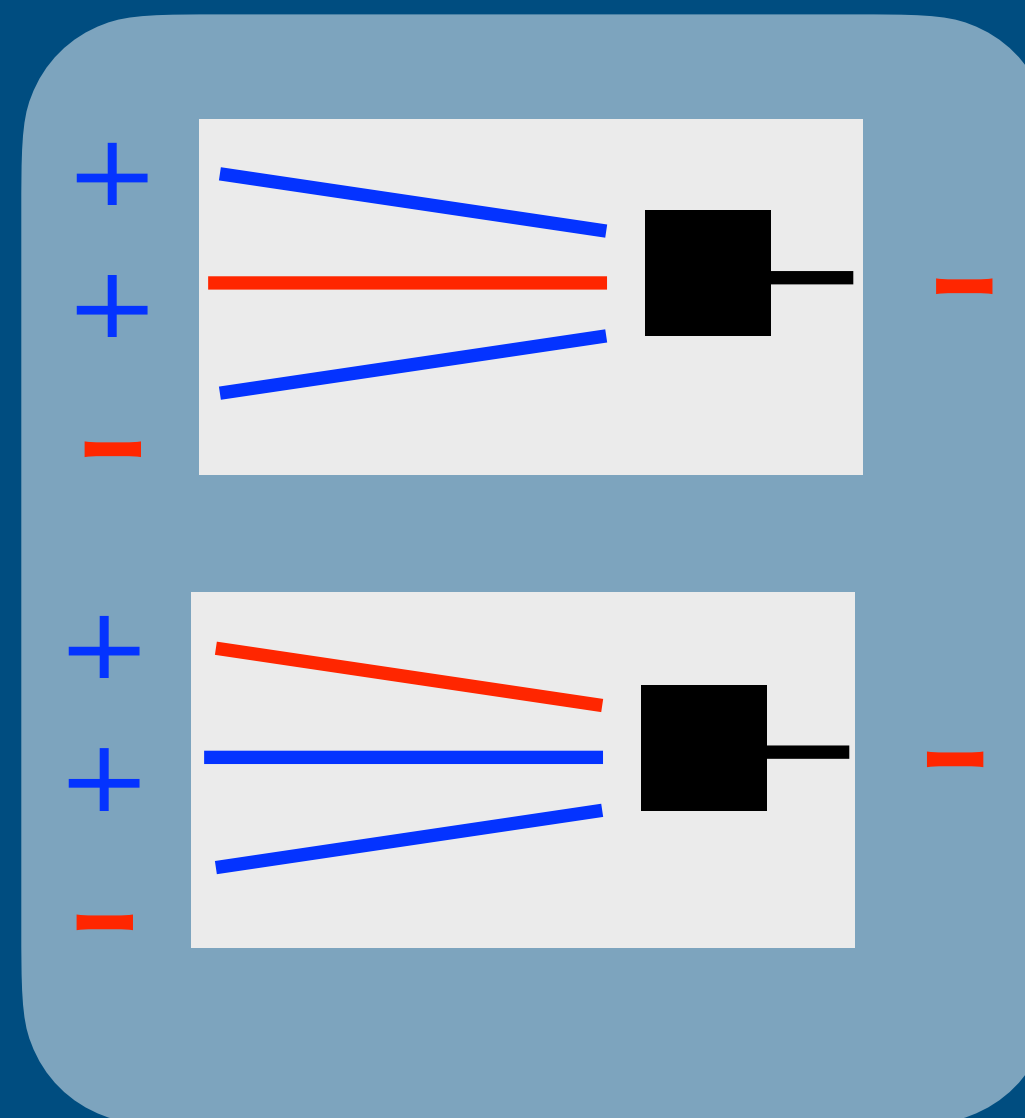
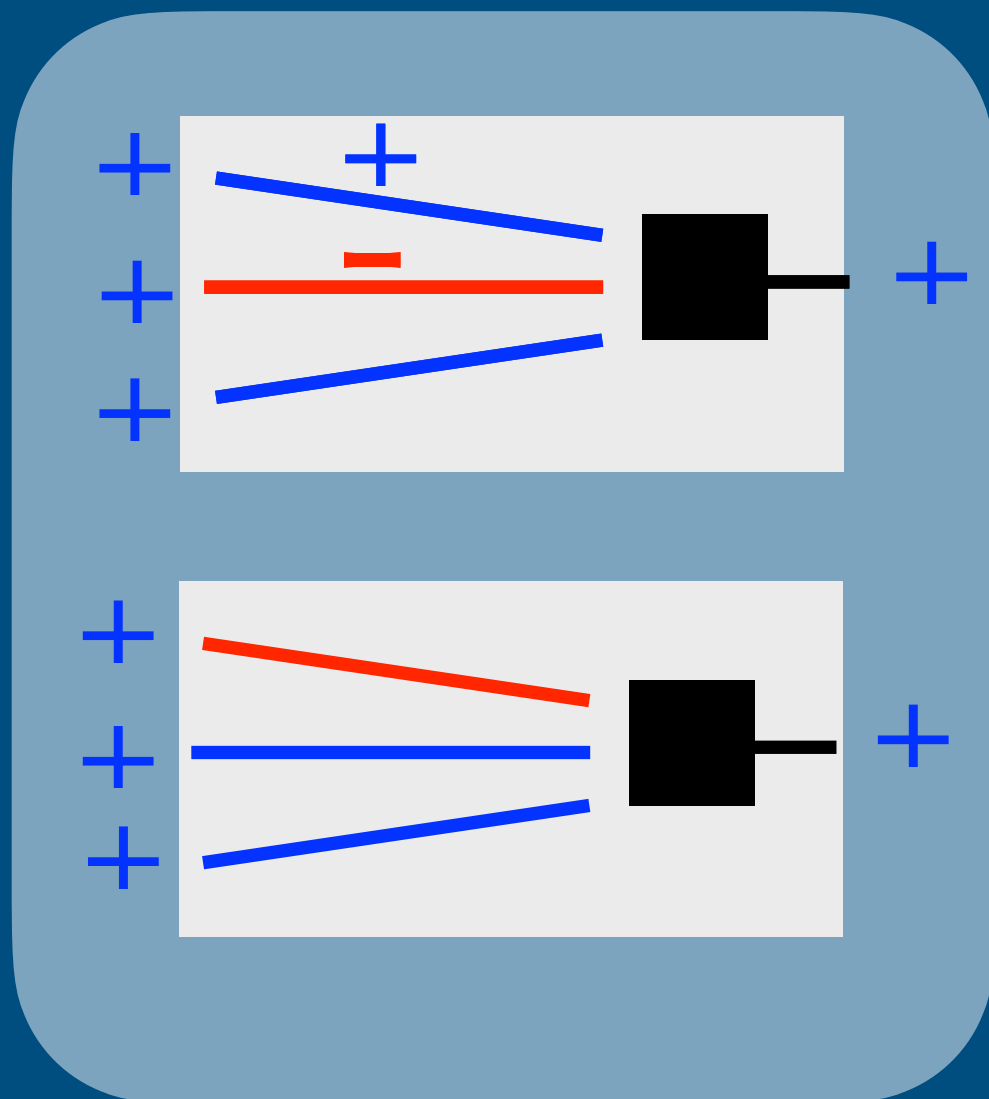
⋮

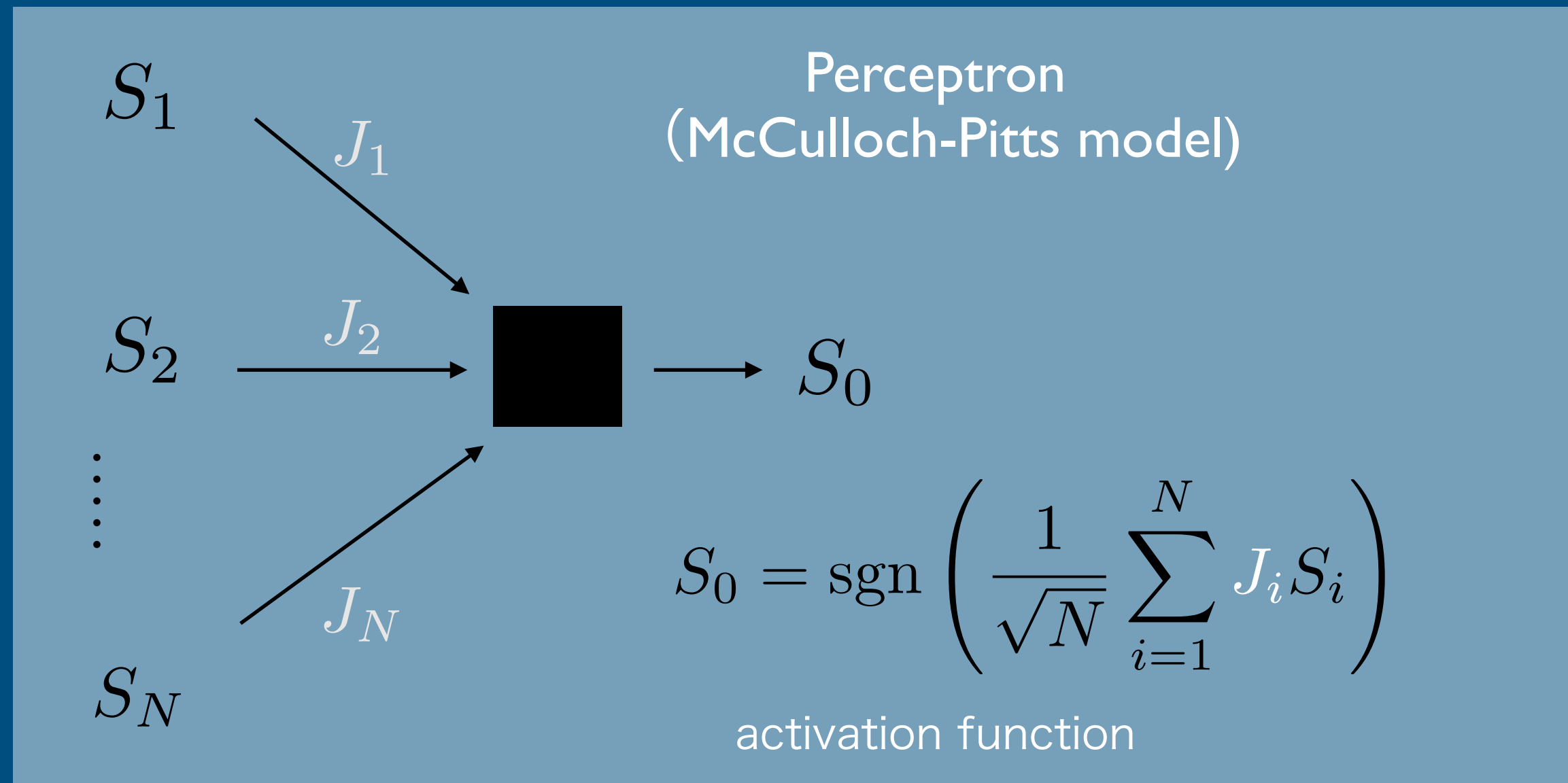
DATA1

DATA2

DATA3

⋯

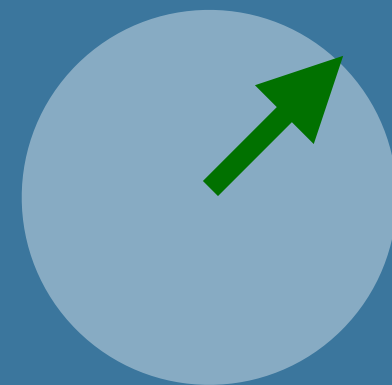




Elisabeth Gardner
(1957-1988)

$$\alpha = \frac{M}{N}$$

$M \rightarrow \infty$ with fixed α



Statistical mechanics of J_i which meet random constraints

$$S_i^\mu = \pm 1 \quad \text{pattern } \mu = 1, 2, \dots, M$$

Gardner volume

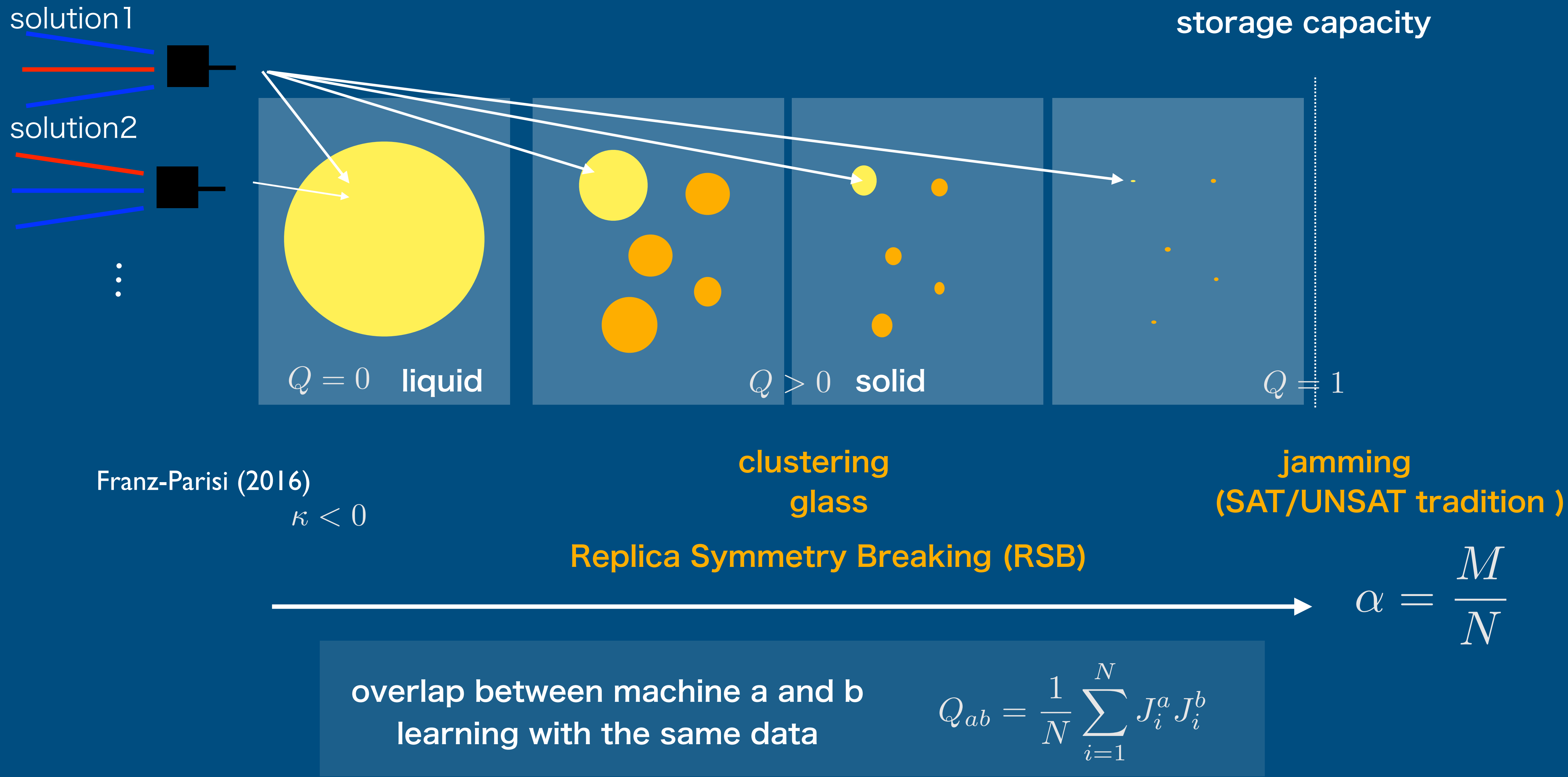
$$V = \int \prod_{j=1}^N \frac{dJ_j}{\sqrt{2\pi}} e^{-\frac{J_j^2}{2}} \prod_{\mu=1}^M e^{-\beta v(r^\mu)}$$

“Hardcore” constraint

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

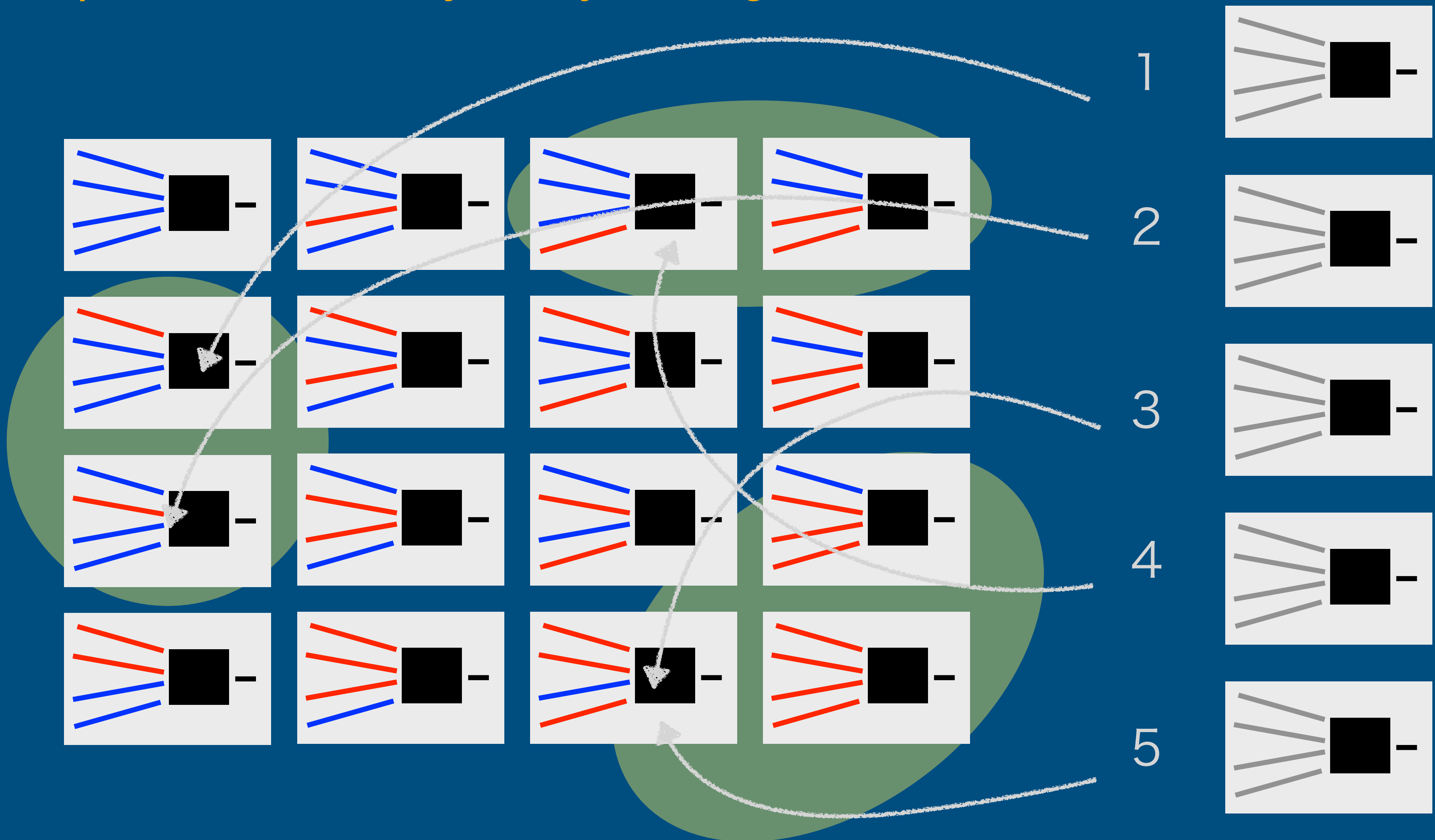
“Gap” $r^\mu = S_0^\mu \sum_{i=1}^N \frac{1}{\sqrt{N}} J_i S_i^\mu - \kappa$

Gardner's volume: design space of a perceptron



**Replica Symmetry Breaking =
Replica "Permutation" Symmetry Breaking**

replica=machines learning in parallel
with the same training data

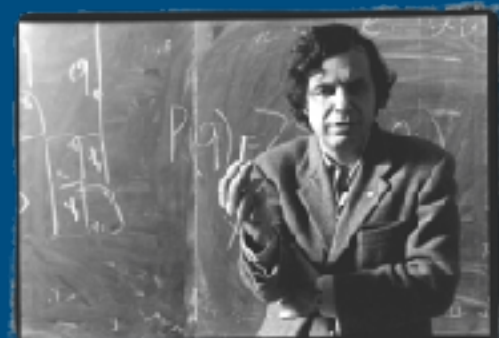
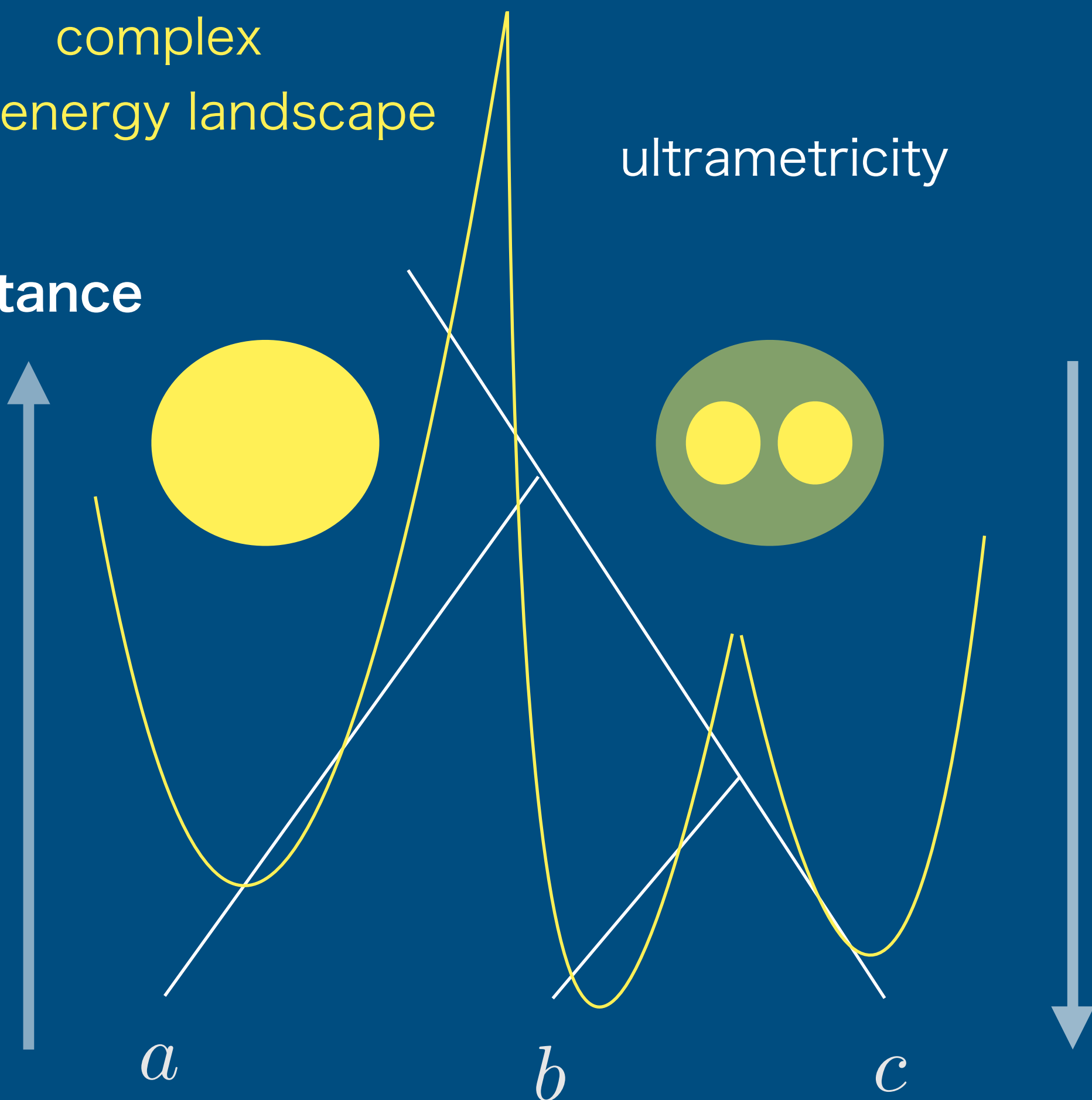


Replica (permutation) symmetry breaking and ultrametricity

complex free-energy landscape

ultrametricity

distance



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

overlap matrix \hat{q}

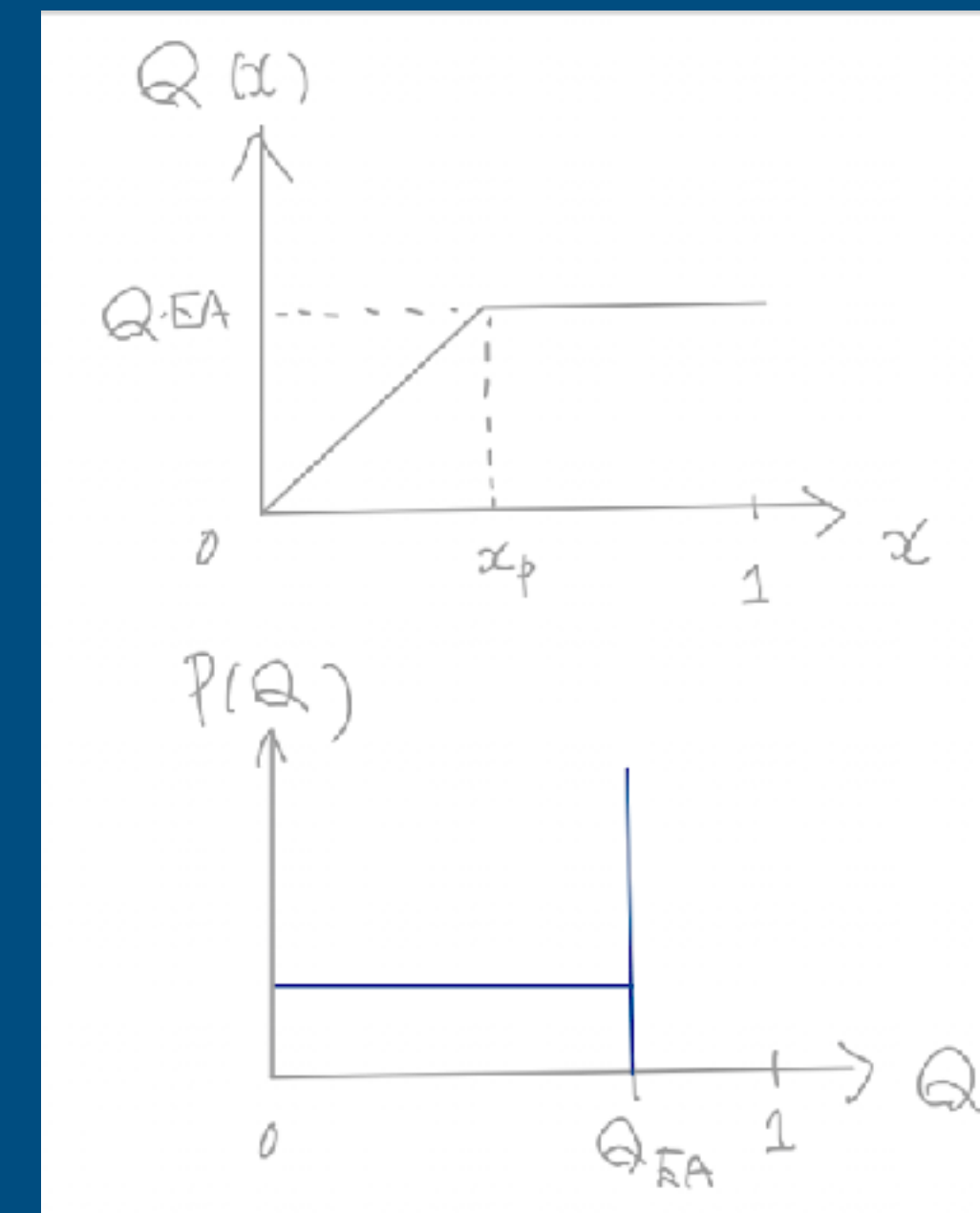
Replica symmetry breaking
G. Parisi (1979)
first found in the SK model
(mean-field spin glass)

Similarity Q

$$Q_{ab} = \frac{1}{N} \sum_{i=1}^N J_i^a J_i^b$$

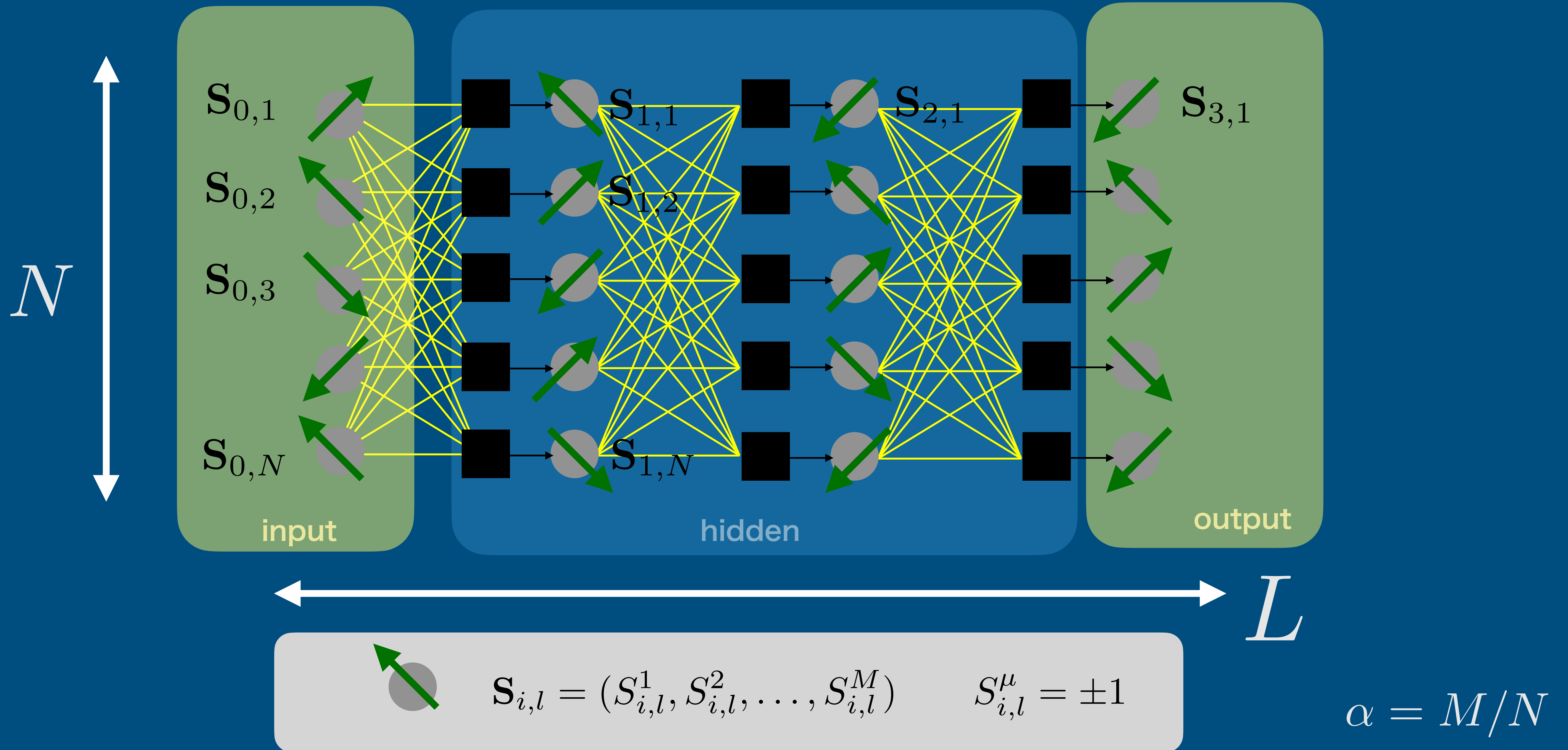
Overlap distribution

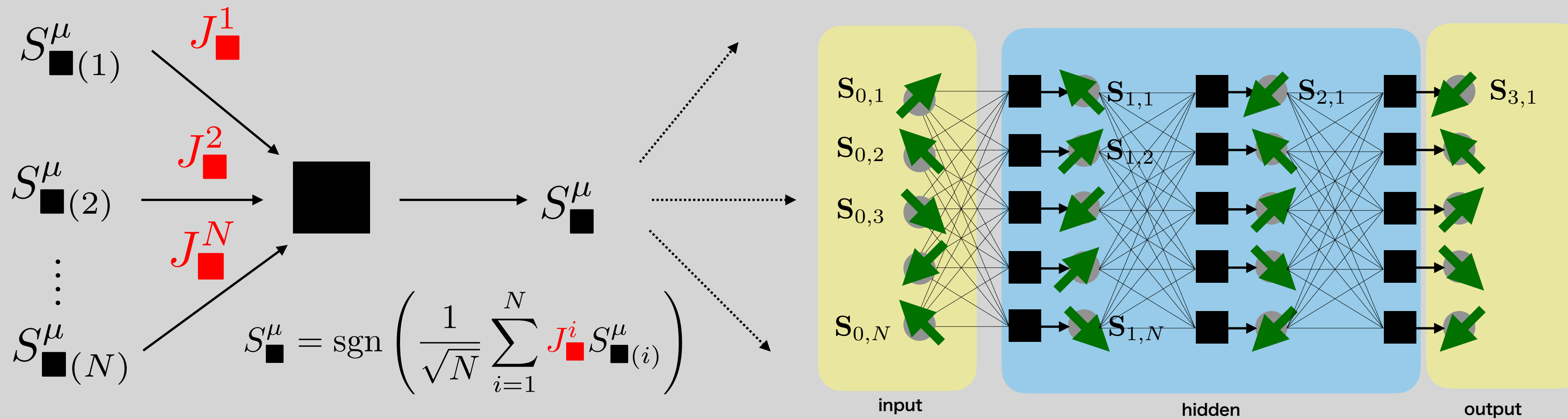
$$P(Q) = \frac{dx}{dQ}$$



Multi-layer Neural Network

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





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

Usual strategy of learning

(1) define "loss function"

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

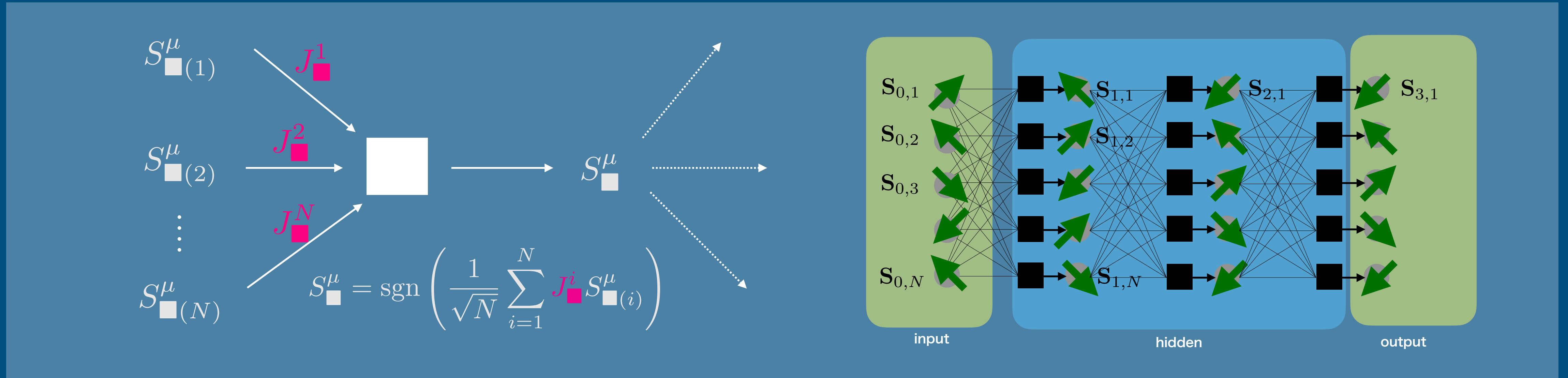
desired output
↓

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

e.g. SDG (stochastic gradient descent)

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

Gardner volume in deep perceptron network



Gardner volume generalized for a multi-layer network

(c.f..) Single perceptron: E. Gardner (1987)

trace over hidden variables

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

Hamiltonian with
“short-ranged” interactions

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

“gap”

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

“Hardcore” constraint

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

We study two scenarios of machine learning with artificially generated training data

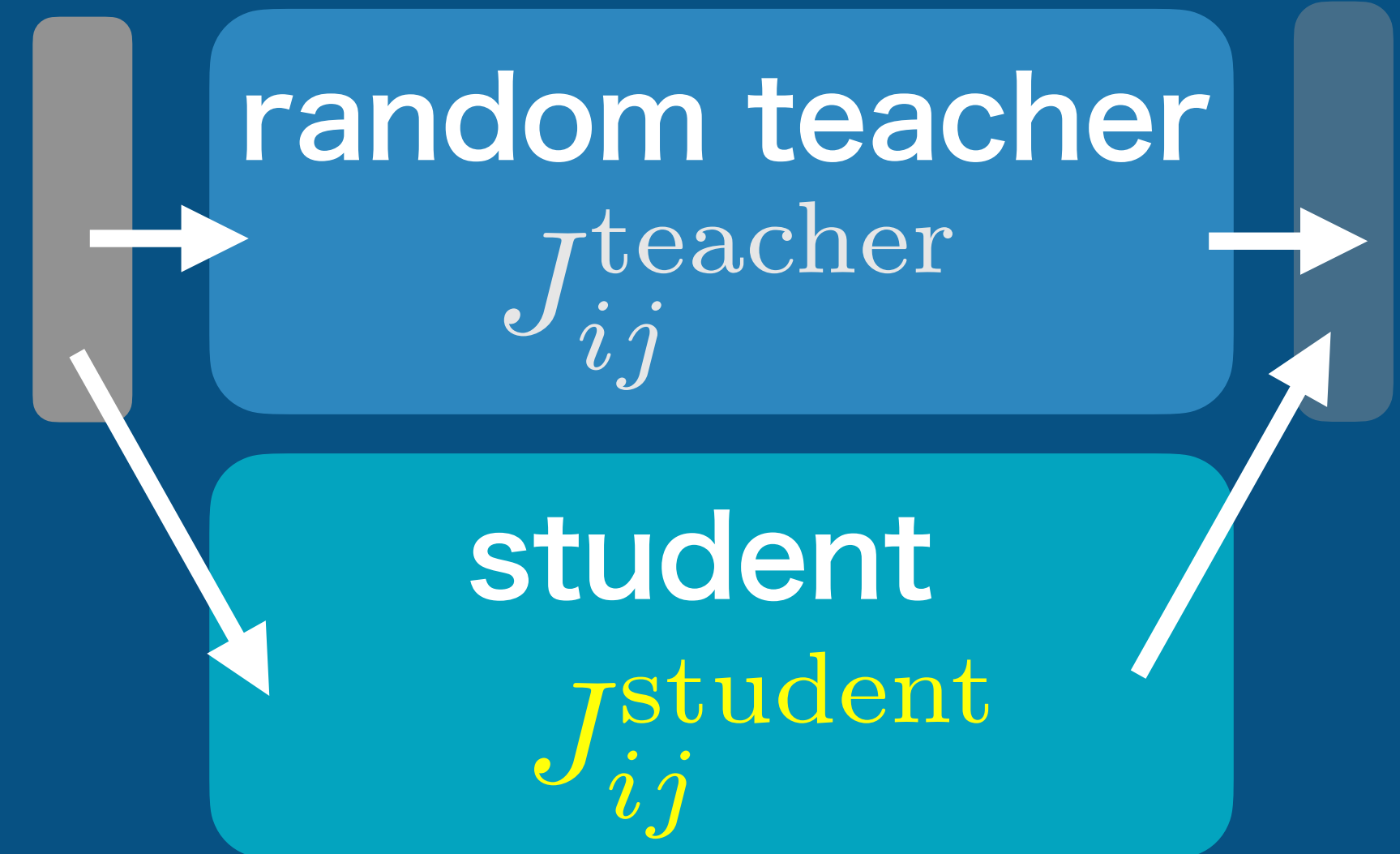
random
input



random
output

scenario (1)

Random
input

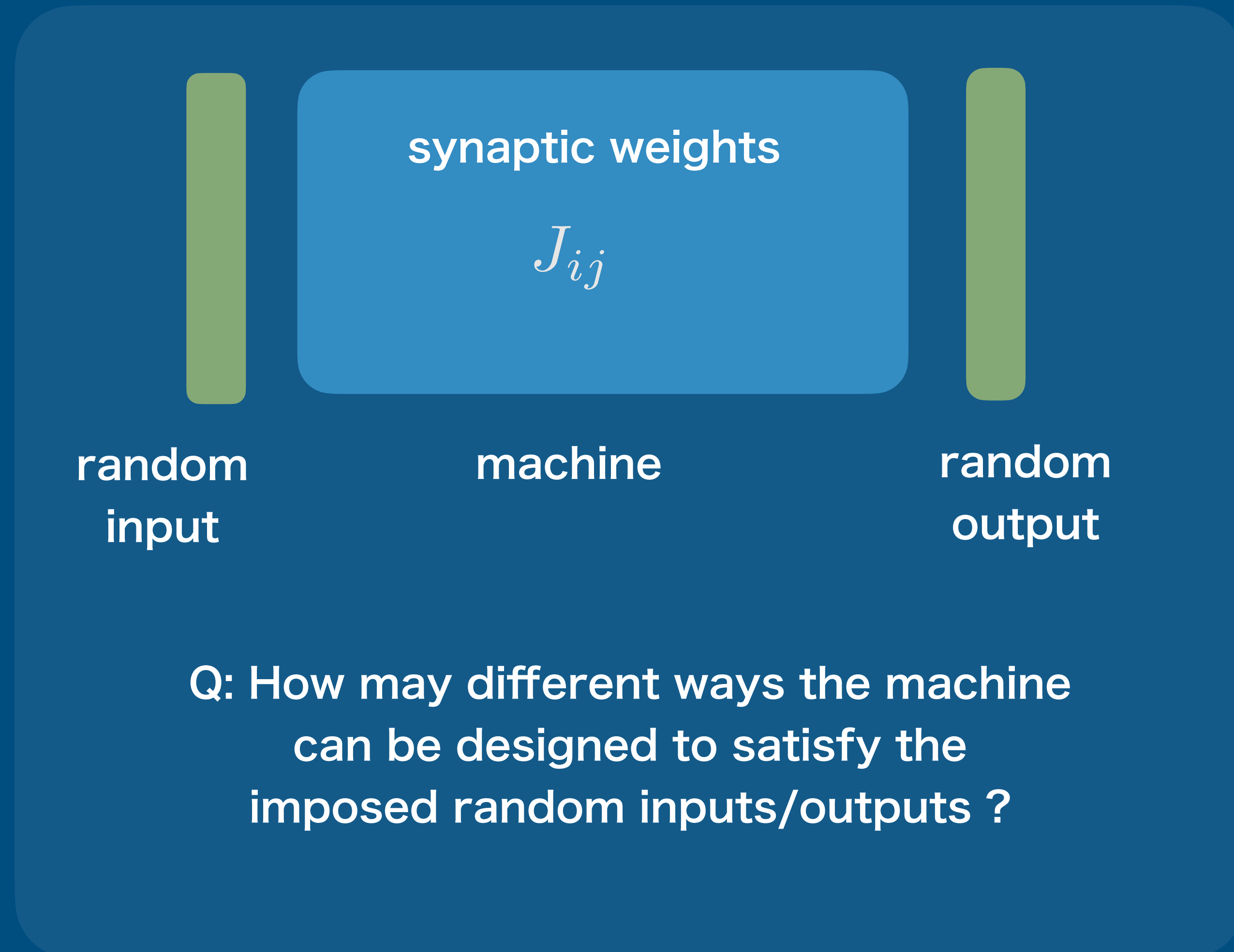


output

scenario (2)

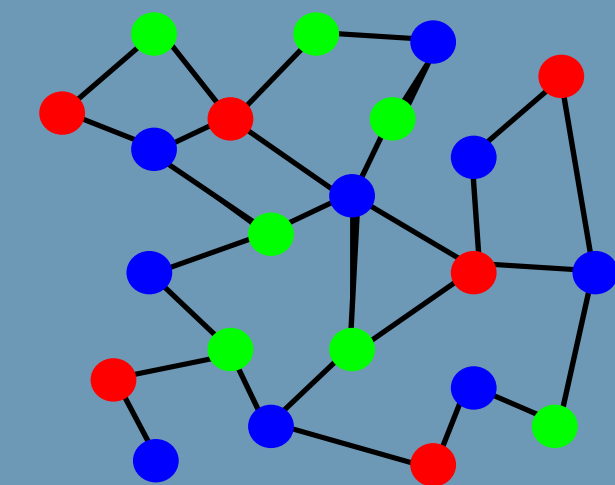
Scenario (I) Random inputs/random outputs

- a constraint satisfaction problem
- glass/jamming physics



“Random Constraint Satisfaction Problem”
(ランダム制約充足問題)

Example:
Graph coloring



Antiferromagnetic Potts model

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

Connection to glass physics

Clustering transition = glass transition
SAT/UNSAT transition = jamming transition

Replicated Gardner volume

replicas: machines learning in parallel with the same data

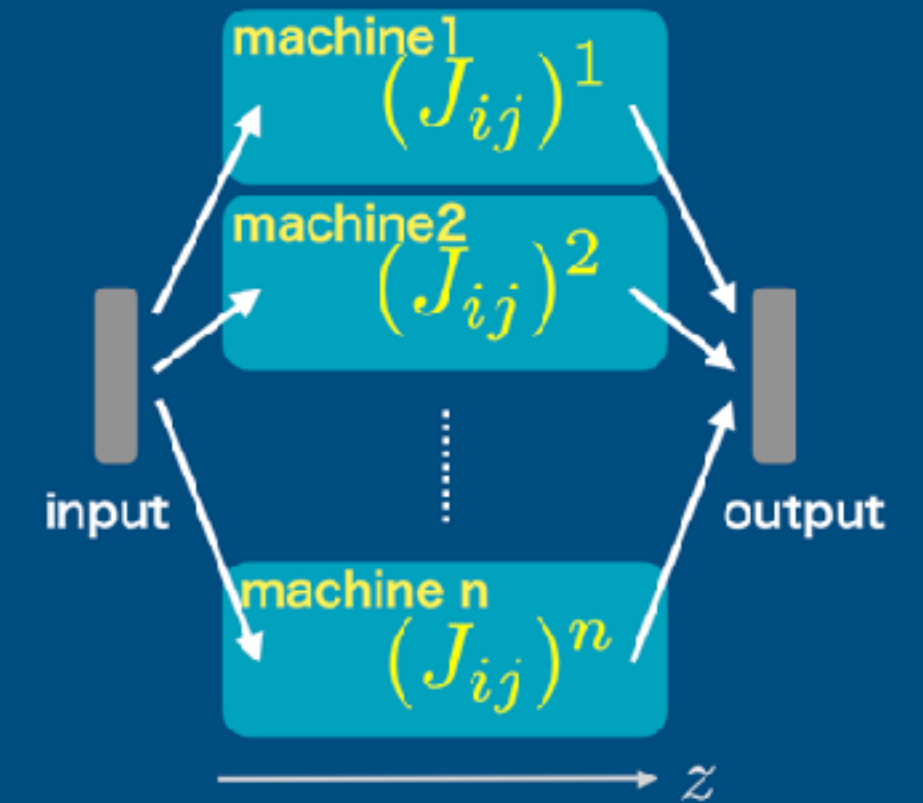
$$V^n(\mathbf{S}_0, \mathbf{S}_L) = \prod_{a=1}^n \left(\prod_{\blacksquare} \text{Tr} \mathbf{J}_{\blacksquare}^a \right) \left(\prod_{\blacksquare \setminus \text{output}} \text{Tr} \mathbf{S}_{\blacksquare}^a \right) \prod_{\mu, \blacksquare, a} e^{-\beta v(r_{\blacksquare, a}^{\mu})}$$

replicated machines $a = 1, 2, \dots, n$

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

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

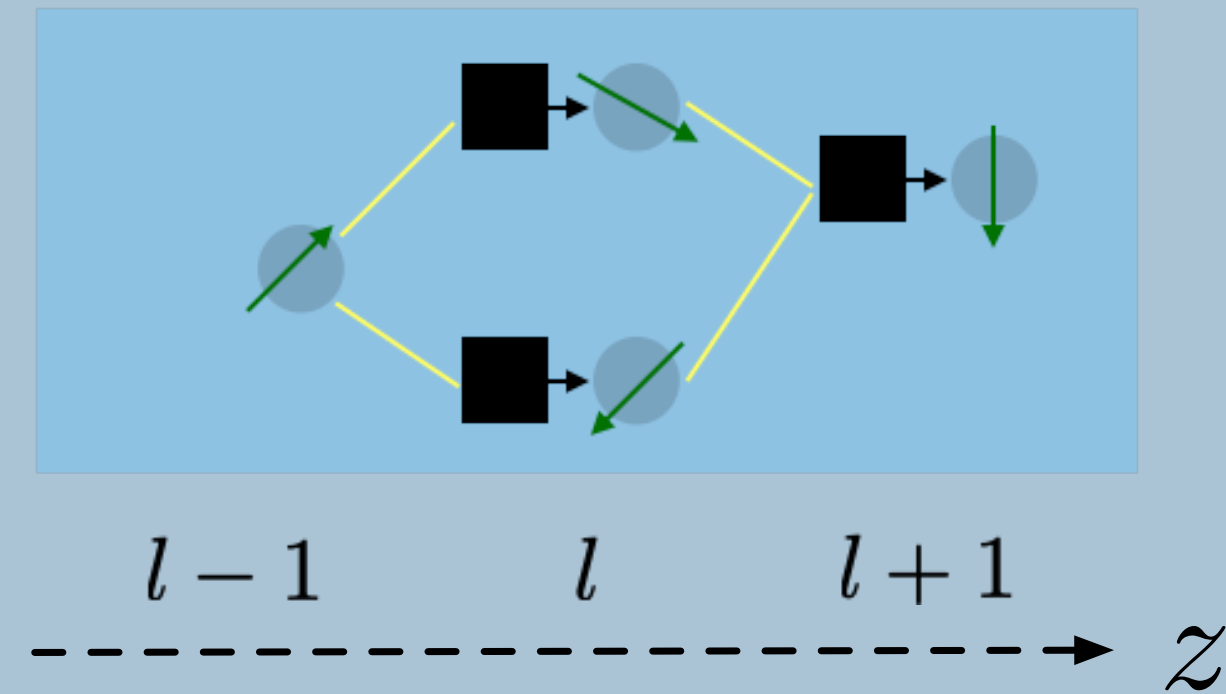
$$\frac{-\beta \overline{F(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} = \frac{\partial_n \overline{V^n(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} \Big|_{n=0} = \partial_n S_n[\{\hat{Q}(l), \hat{q}(l)\}] \Big|_{n=0}$$

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

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

$\alpha = \frac{M}{N}$

Loop correction

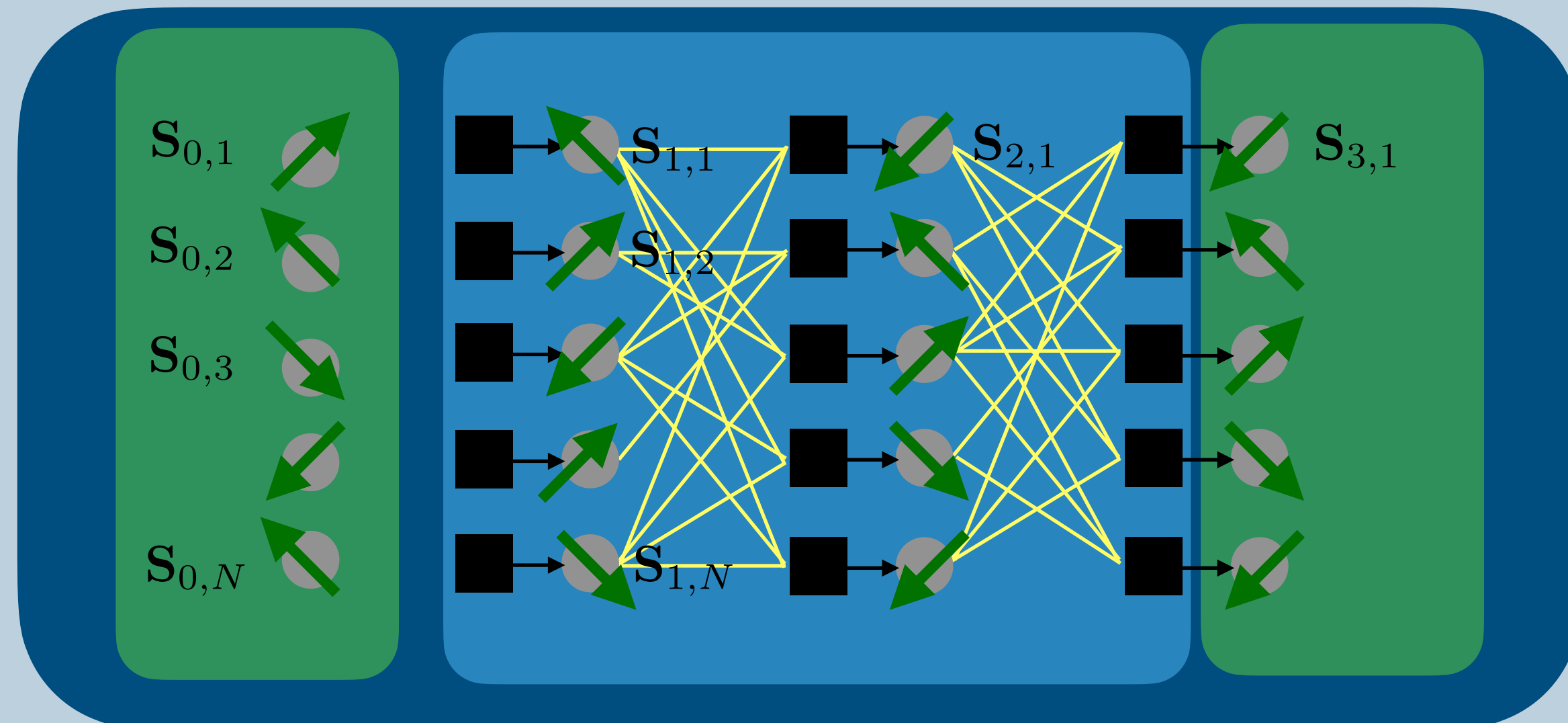


Not global but **dense** coupling - loop corrections become negligible

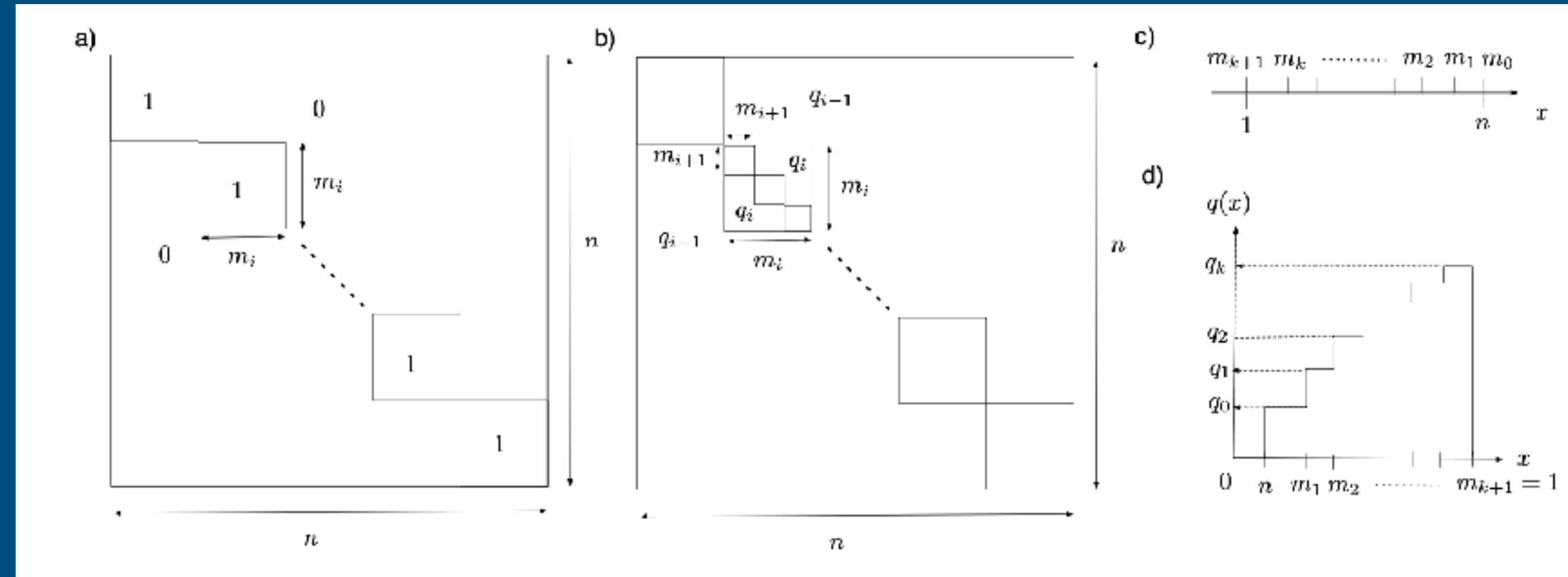
“connectivity” c

$$N \gg c \gg 1$$

$$\alpha = M/c$$



■ Parisi's RSB ansatz



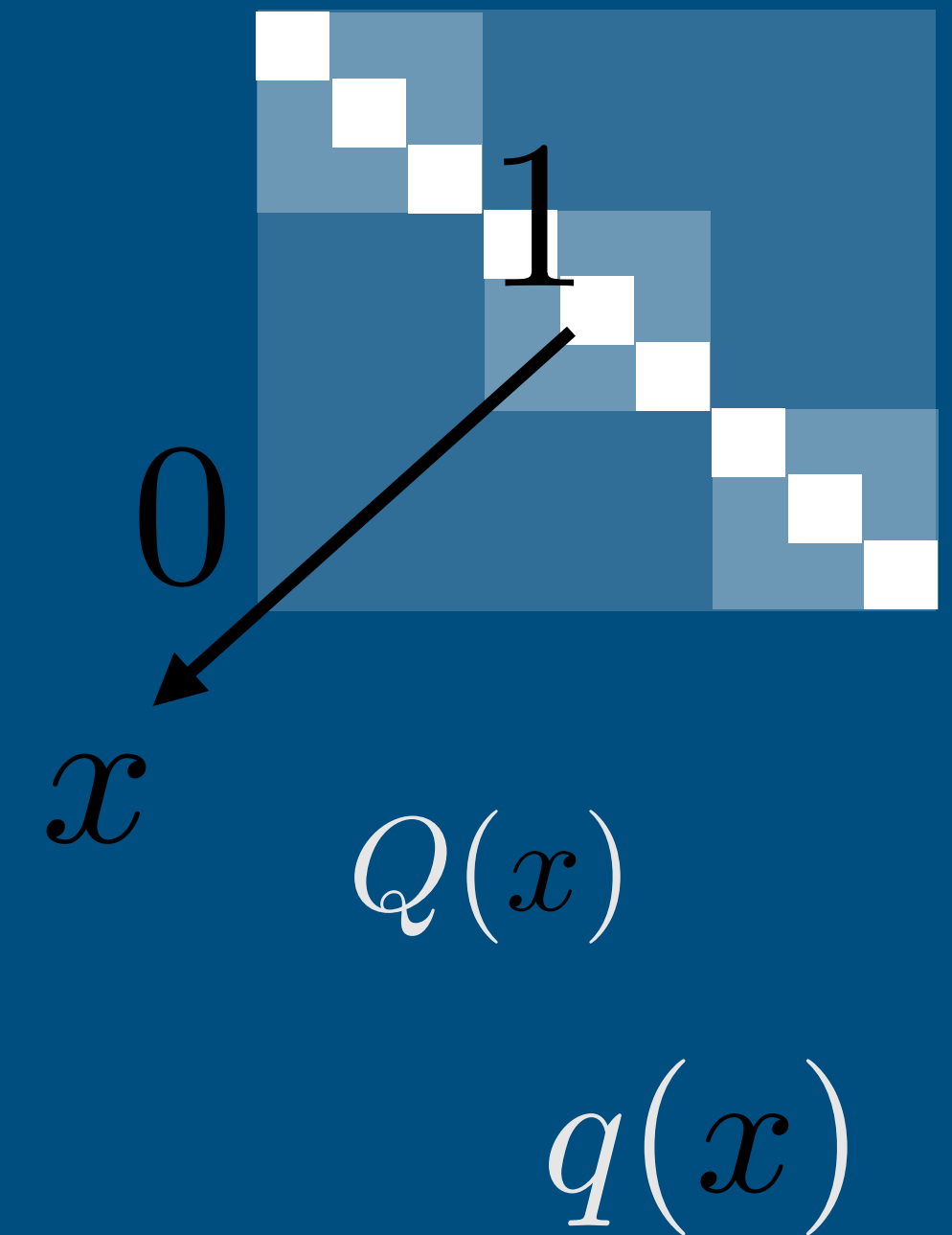
$$Q_{ab}(l) = \sum_{i=0}^{k+1} Q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L$$

$$q_{ab}(l) = \sum_{i=0}^{k+1} q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L-1$$

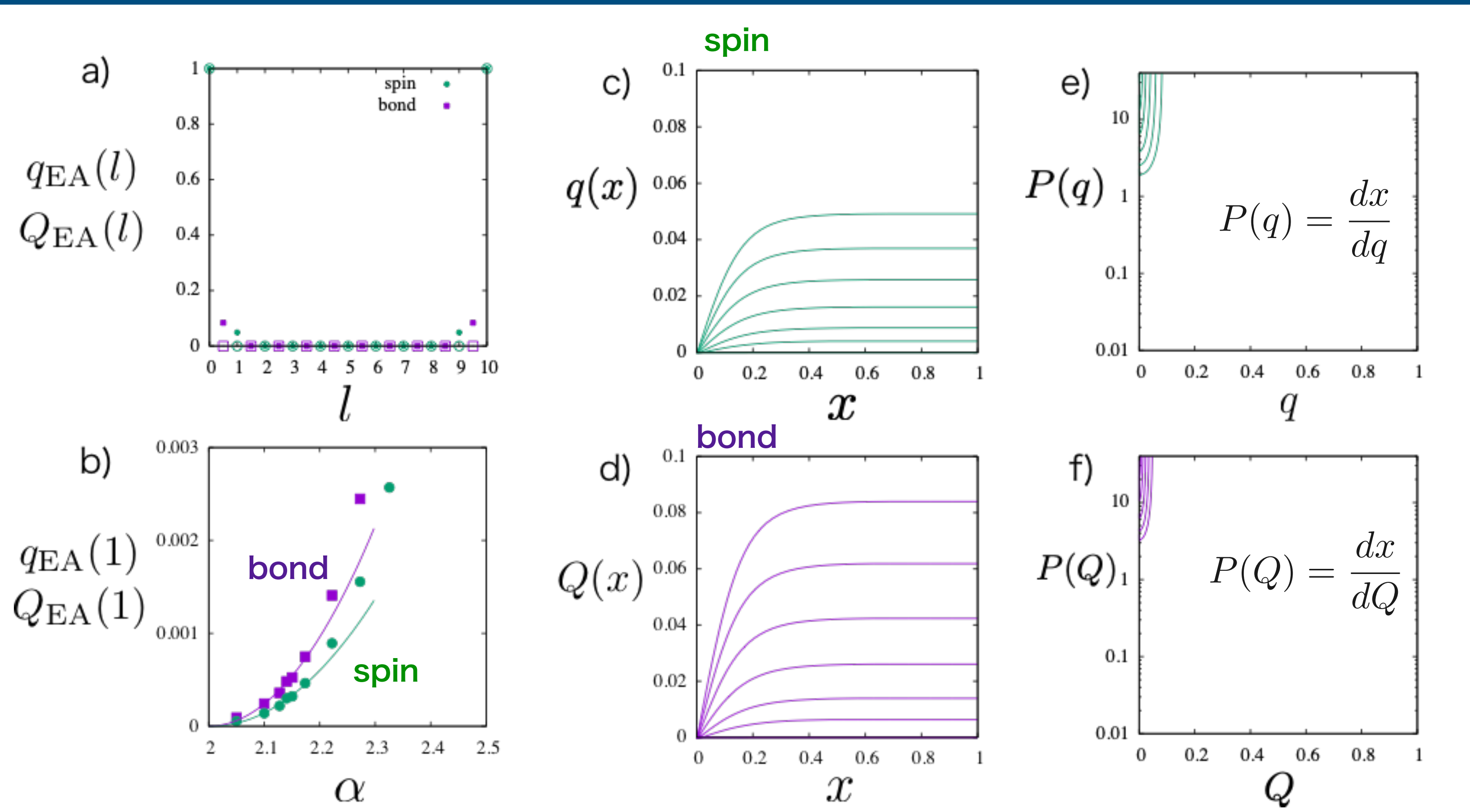
■ Input/output boundaries

replicated machines are subjected to the same training data

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



1st Glass transition



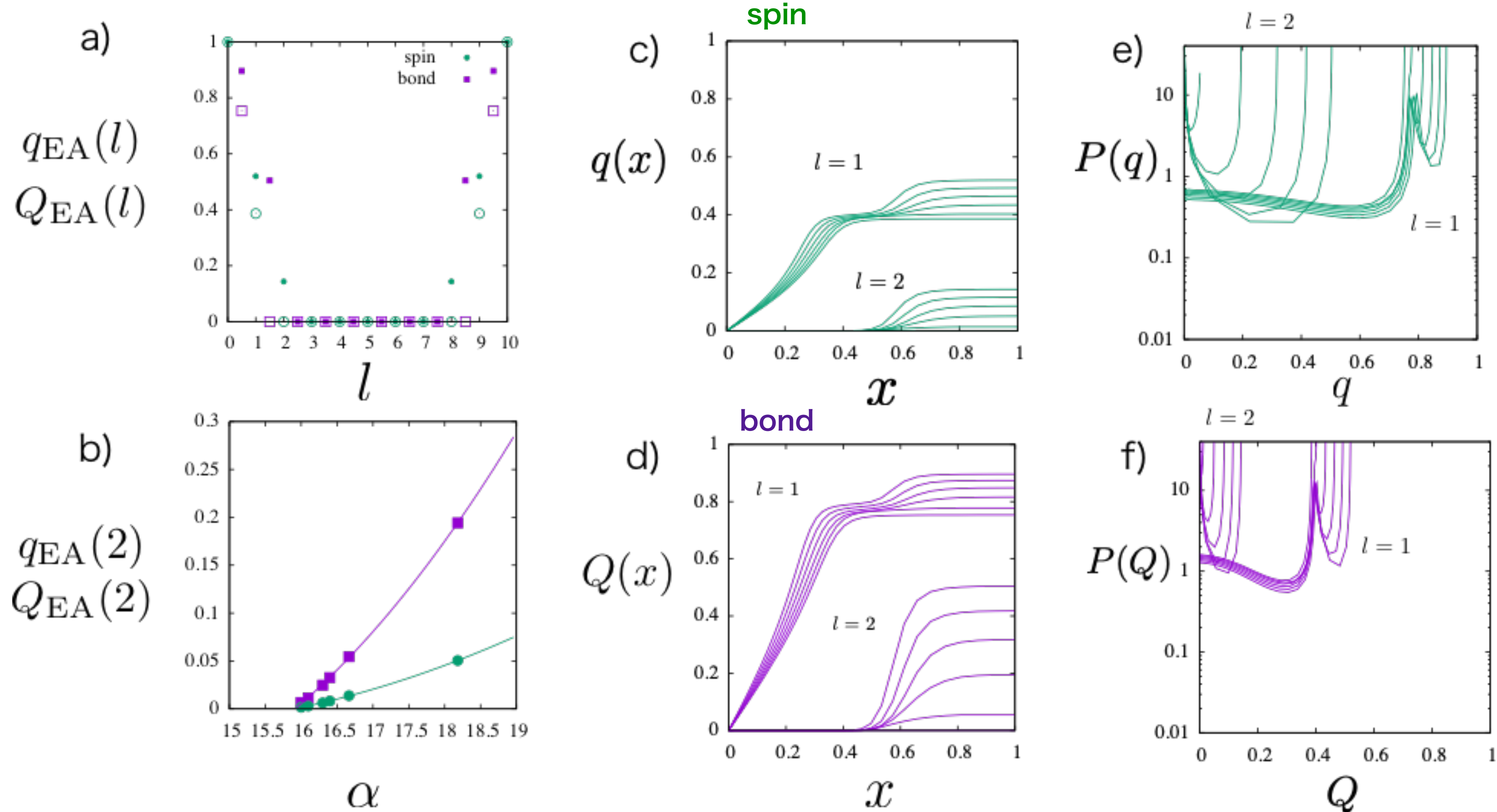
$$\alpha_g \simeq 2.03$$

continuous transition to full RSB glass phase

at 1st & (L-1)th layer

other layers remain in the liquid phase

■ 2nd Glass transition



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

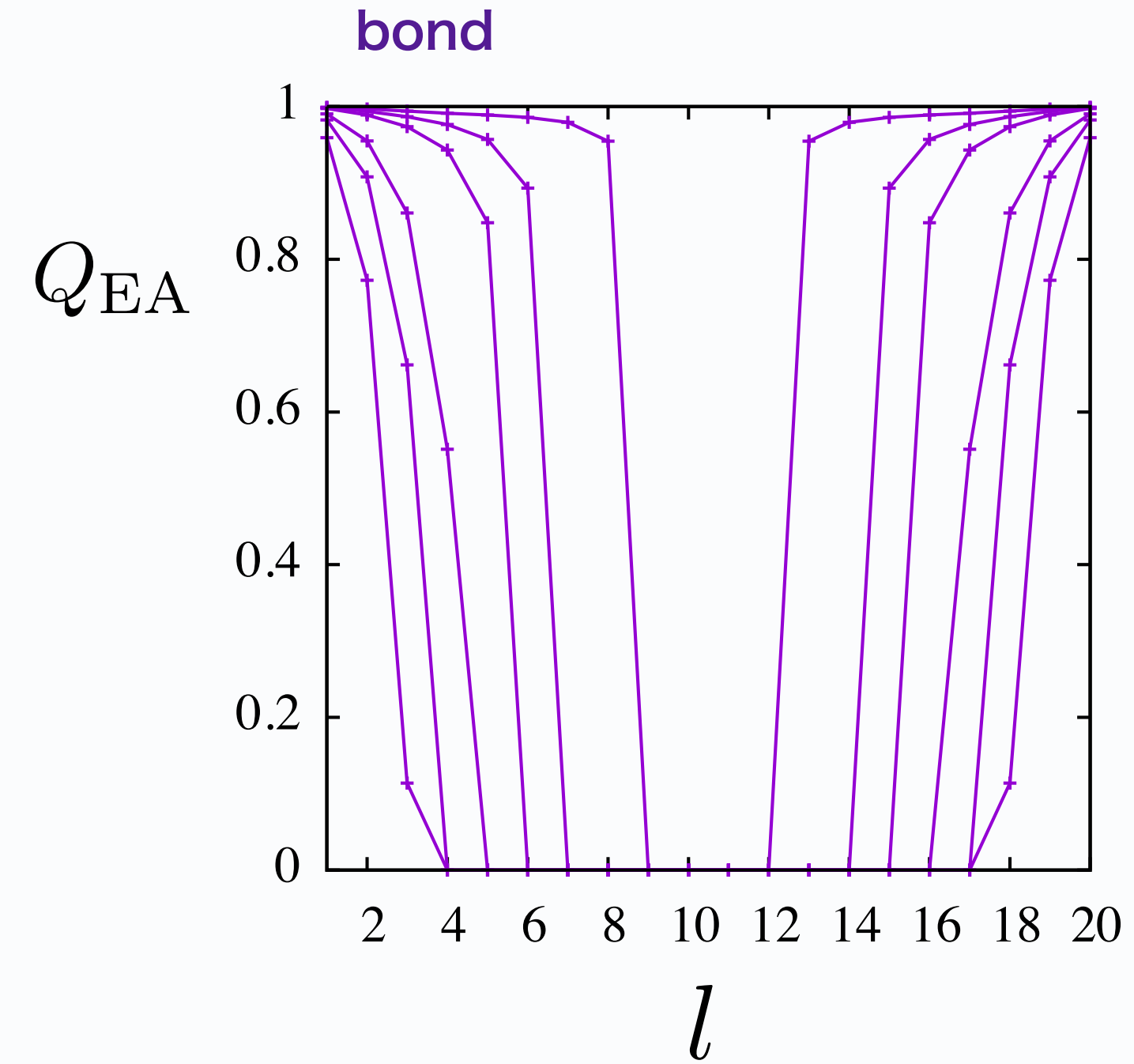
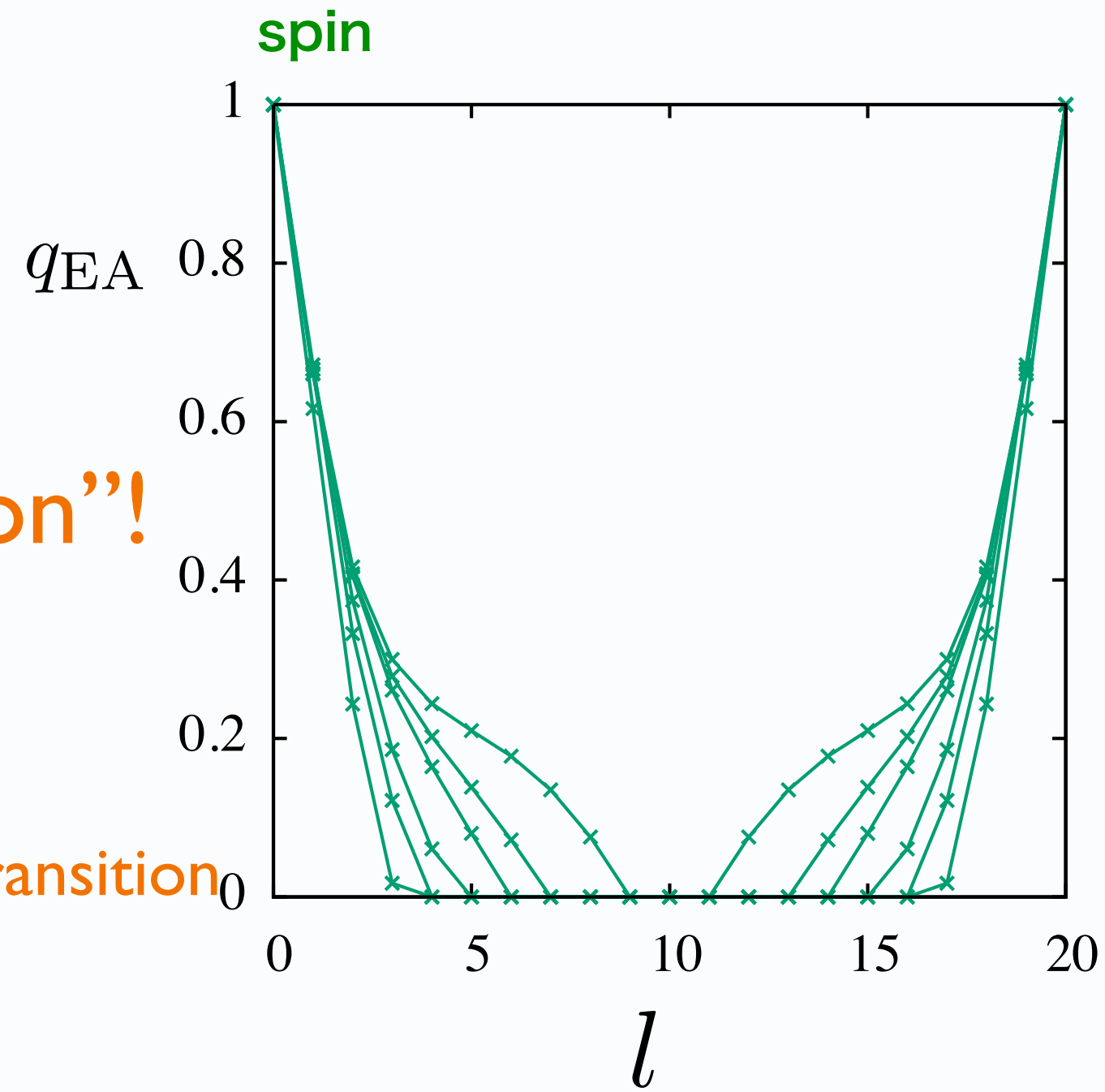
continuous transition to full RSB glass phase
at 2nd & (L-2) th layer

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

■ Growth of glass phase with increasing training data

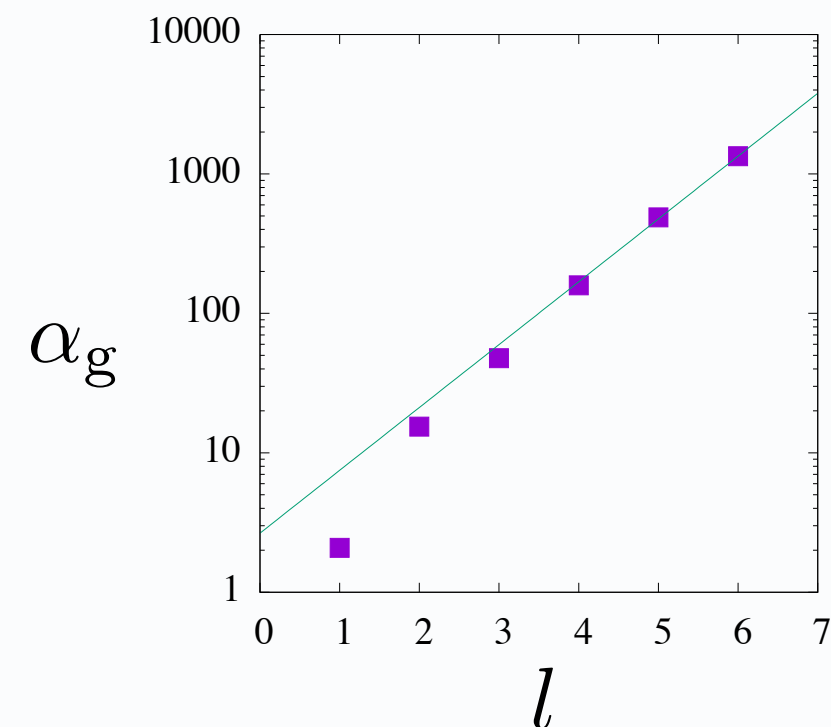
“Wetting transition”!

layer-by-layer 2nd order transition



$\alpha = 50, 100, 200, 1000, 2000, 4000$

“penetration depth”



$$\alpha_g(l) \sim 2.7(3)e^{1.03(2)l}$$

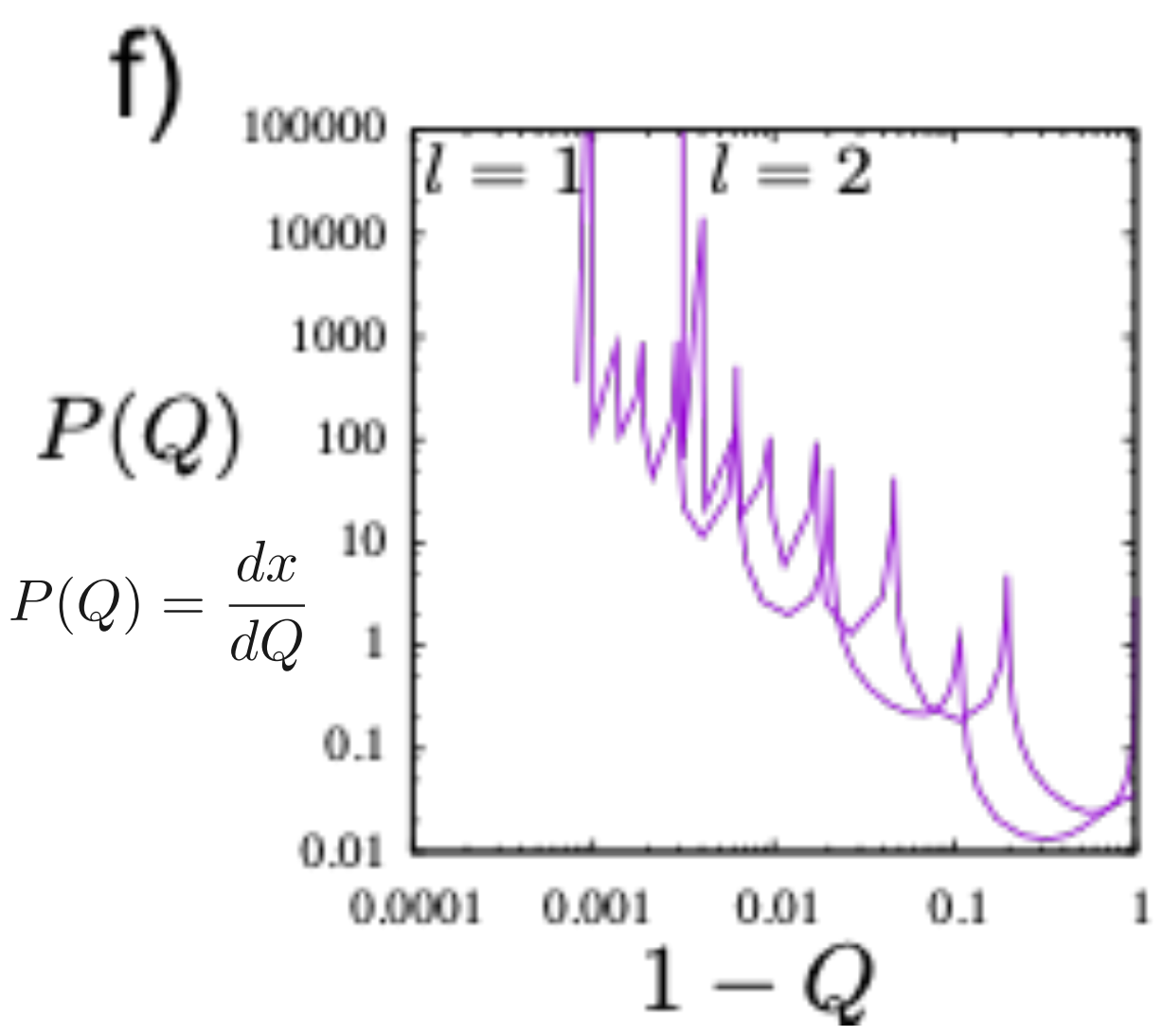
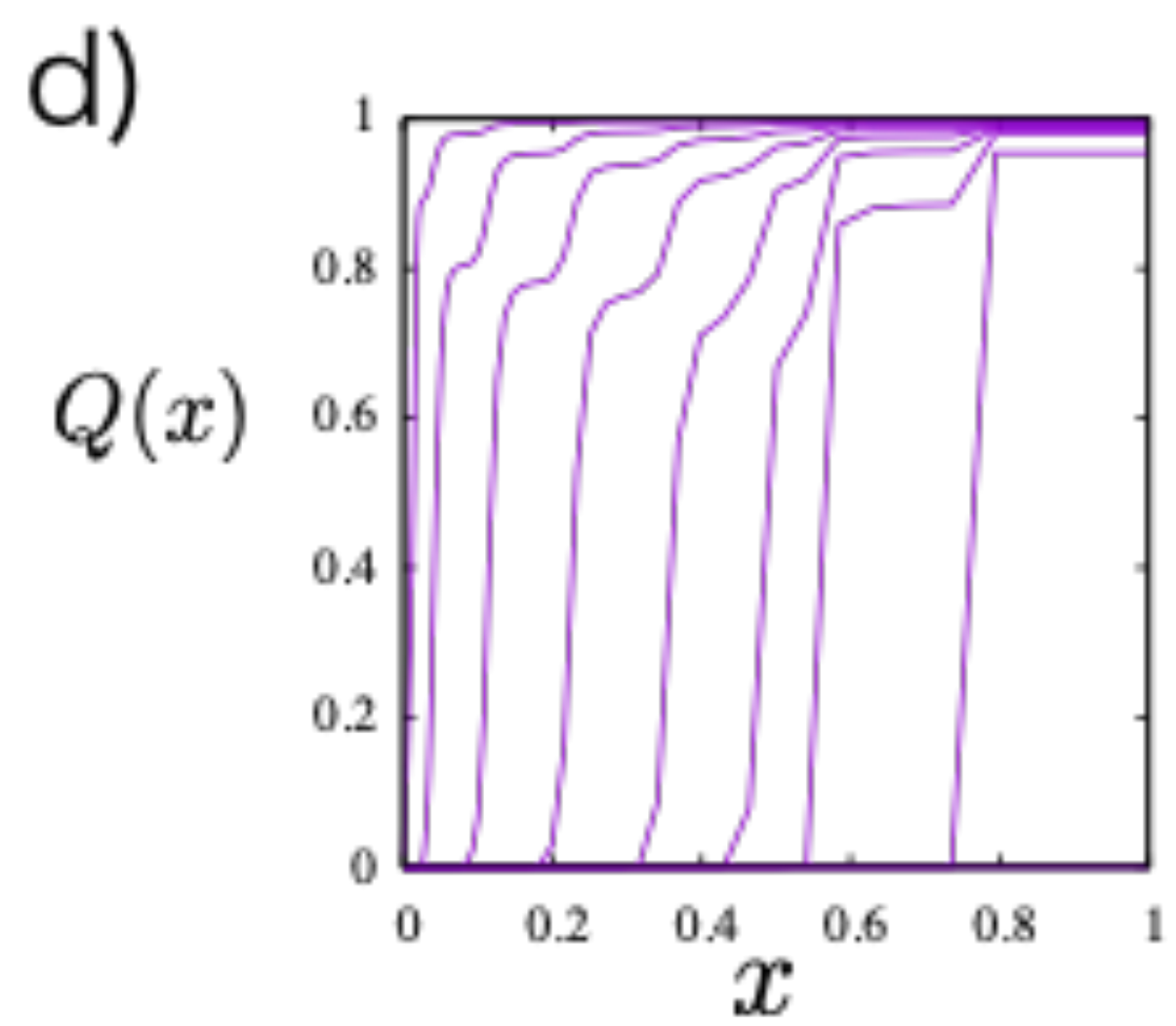
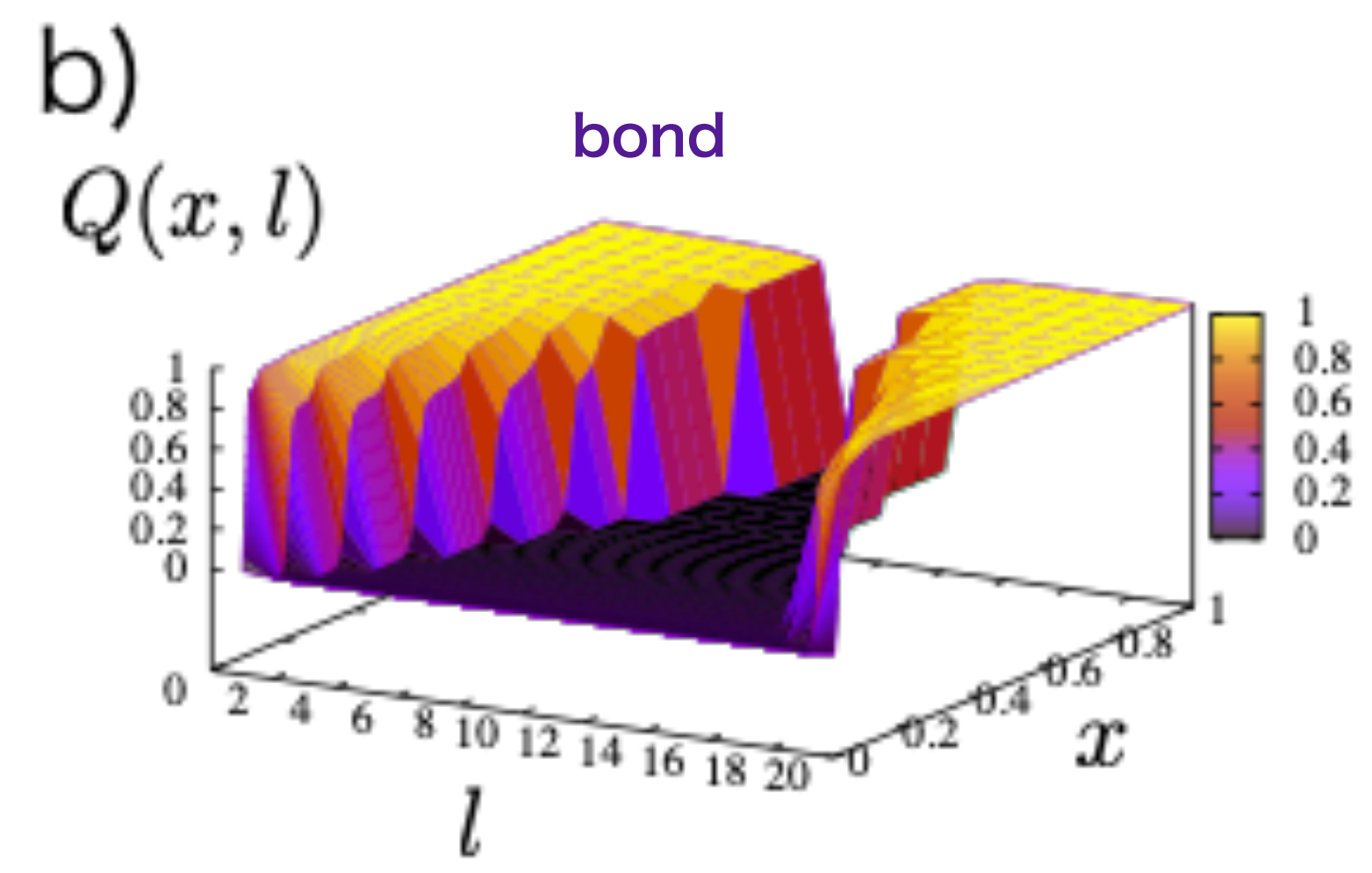
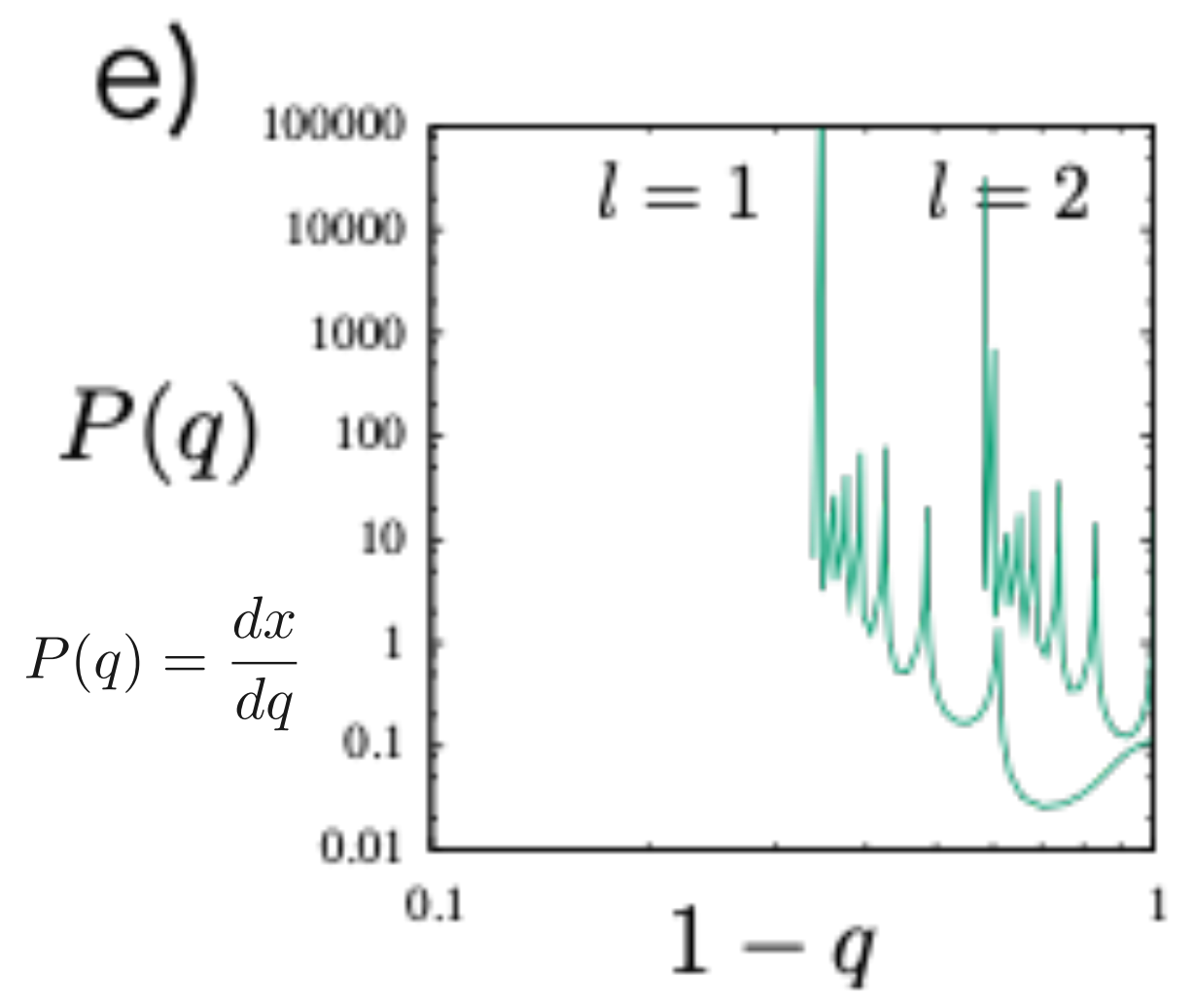
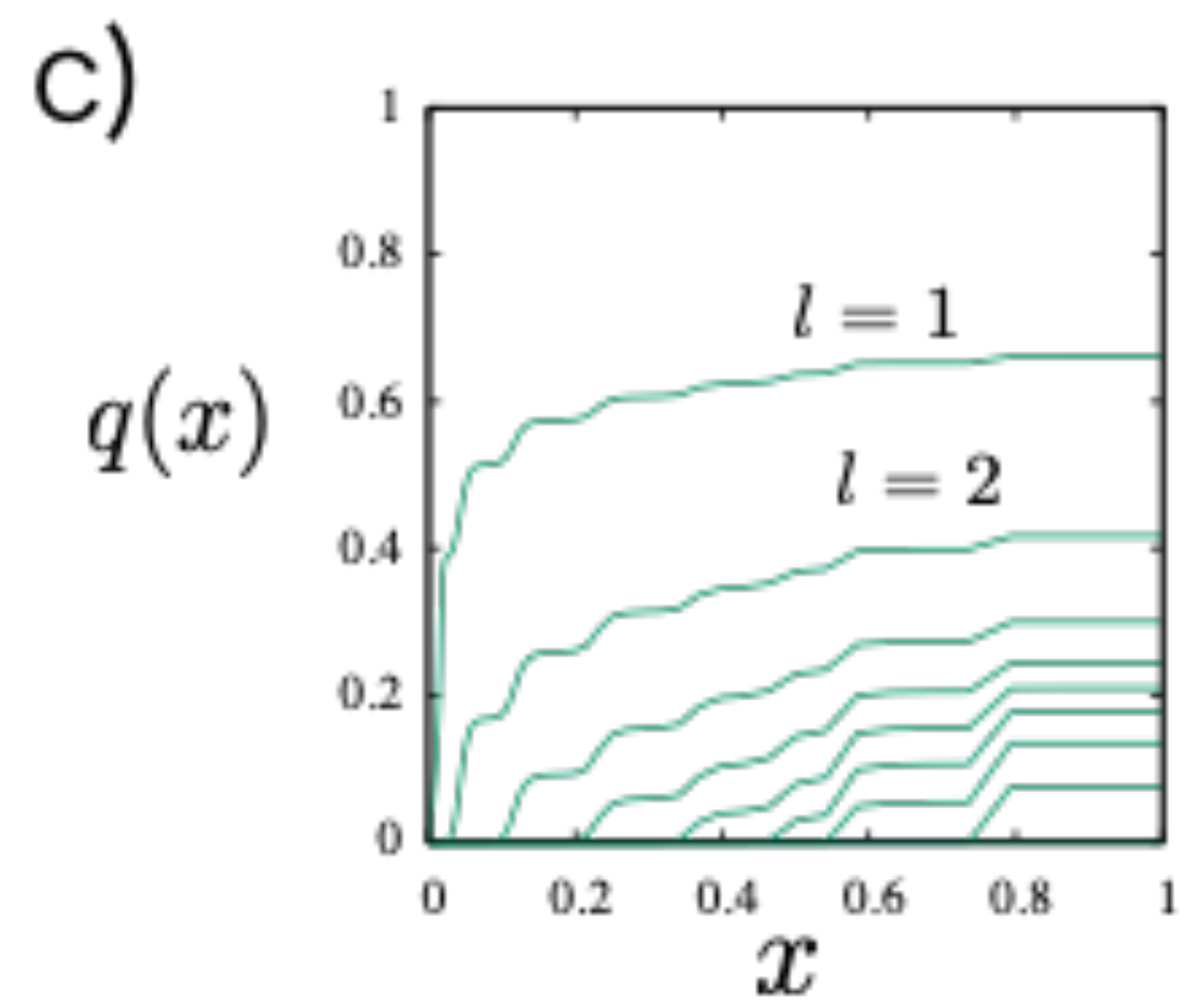
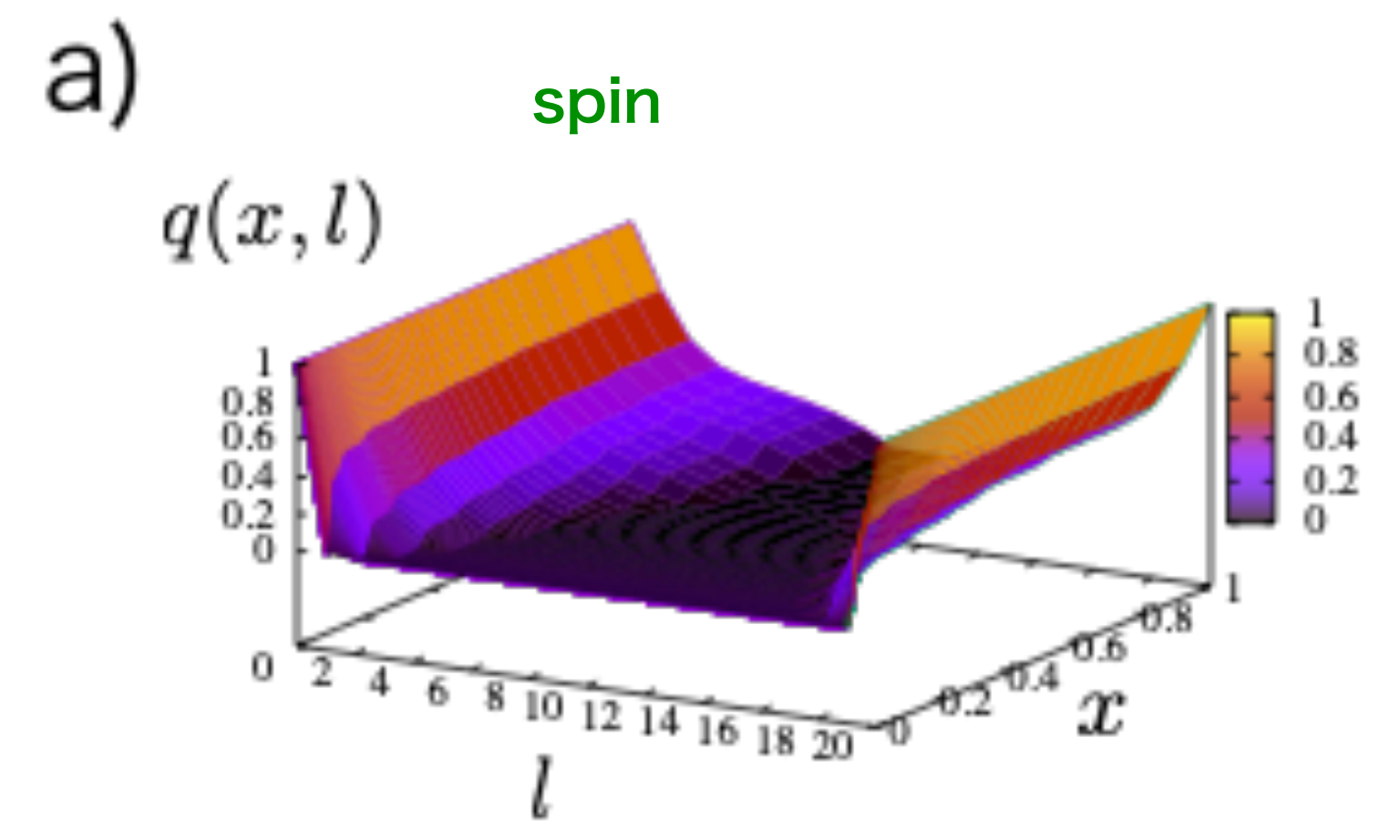
$$\xi(\alpha) \propto \ln \alpha$$

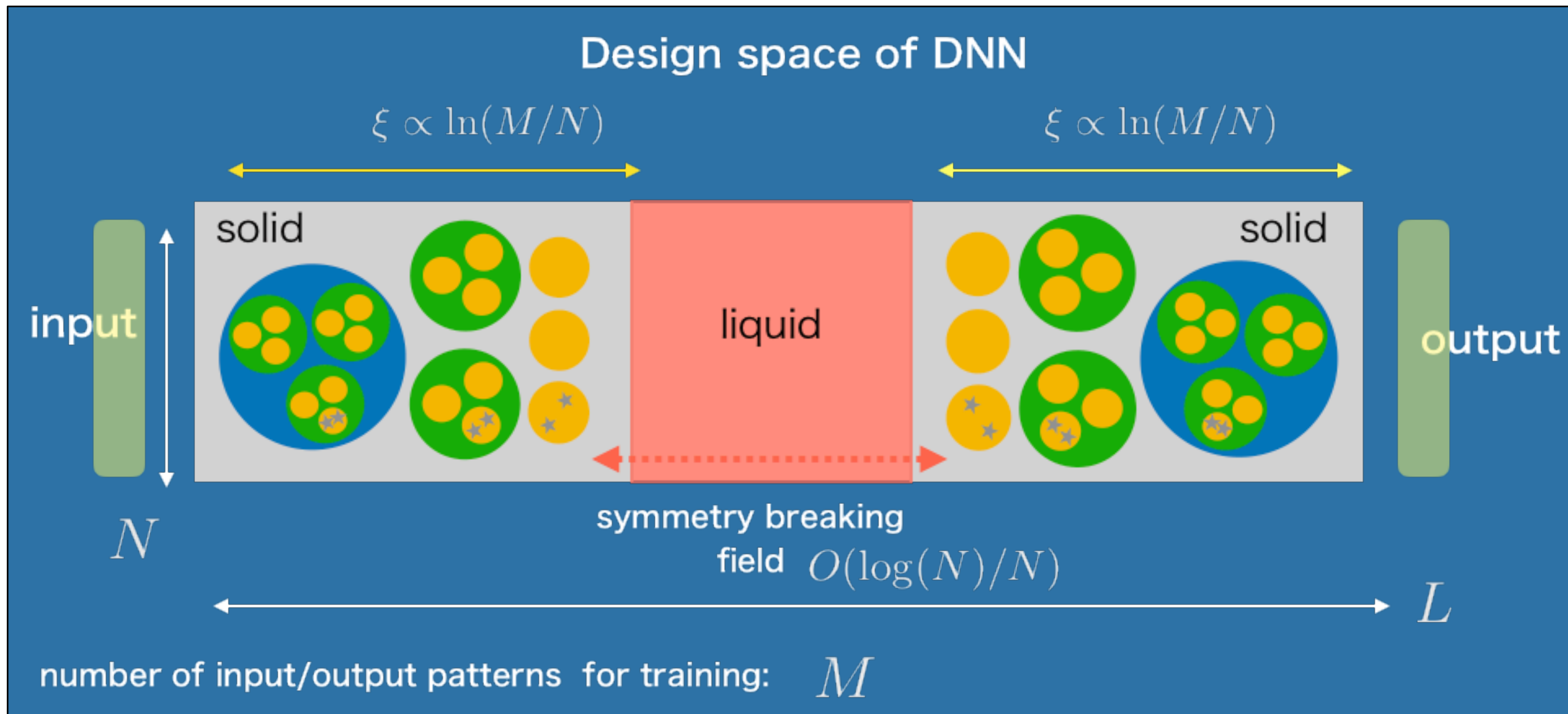
this suggests storage capacity also grows fast with the depth

$$\alpha_J(l) \propto e^{\text{const}l}$$

Space-dependent replica-symmetry breaking

$\alpha = 4000$

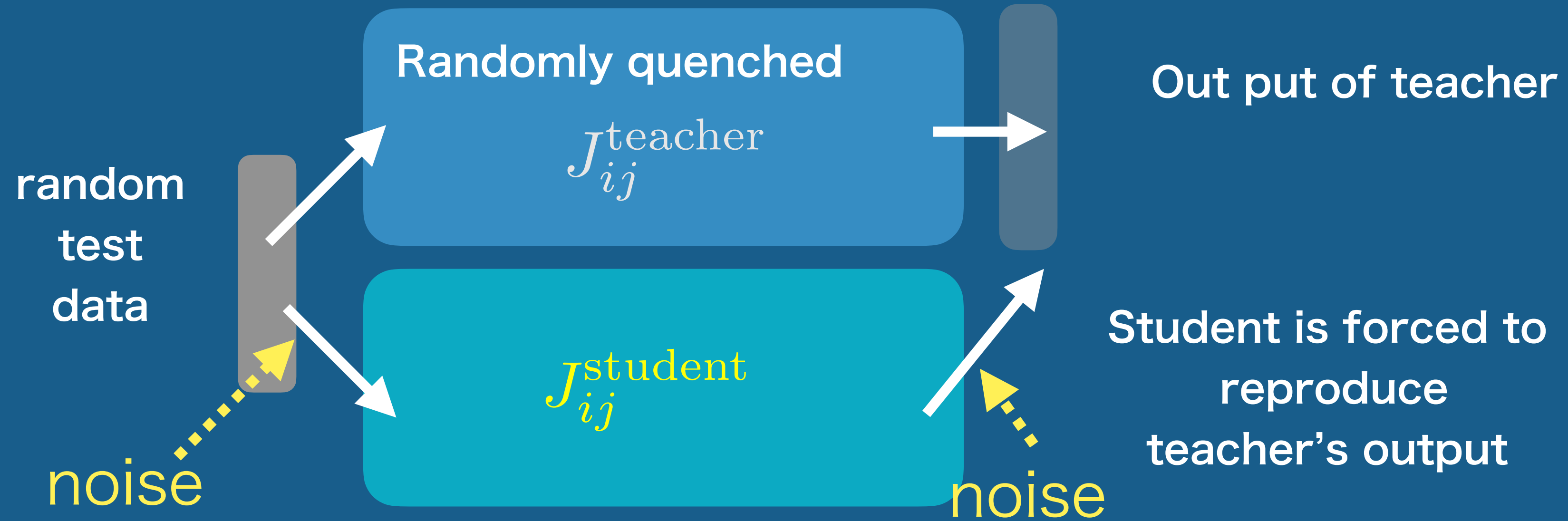




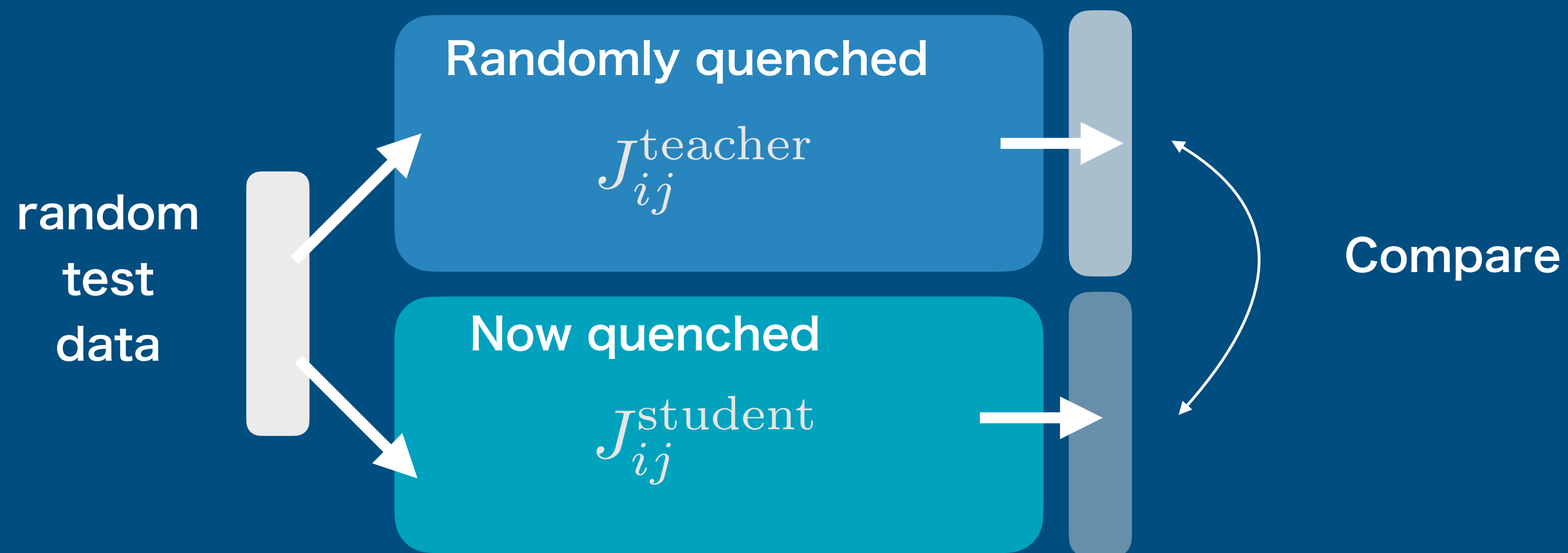
each DNN is NOT a glass

Scenario (2) Teacher student setting - a statistical inference problem

1) Training



2) Test



Replicated Gardner volume

$$V^{1+n}(\mathbf{S}_0, \mathbf{S}_L) = \prod_{a=0}^n \left(\prod_{\blacksquare} \text{Tr} J_{\blacksquare}^a \right) \left(\prod_{\blacksquare \setminus \text{output}} \text{Tr} S_{\blacksquare}^a \right) \prod_{\mu, \blacksquare, a} e^{-\beta v(r_{\blacksquare, a}^{\mu})} \quad r_{\blacksquare, a}^{\mu} = S_{\blacksquare, a}^{\mu} \sum_{i=1}^N \frac{1}{\sqrt{N}} J_{\blacksquare, a}^i S_{\blacksquare(i), a}^{\mu}$$

teacher-machine $a = 0$ **student-machines** $a = 1, 2, \dots, n$

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
(Franz-Parisi potential)

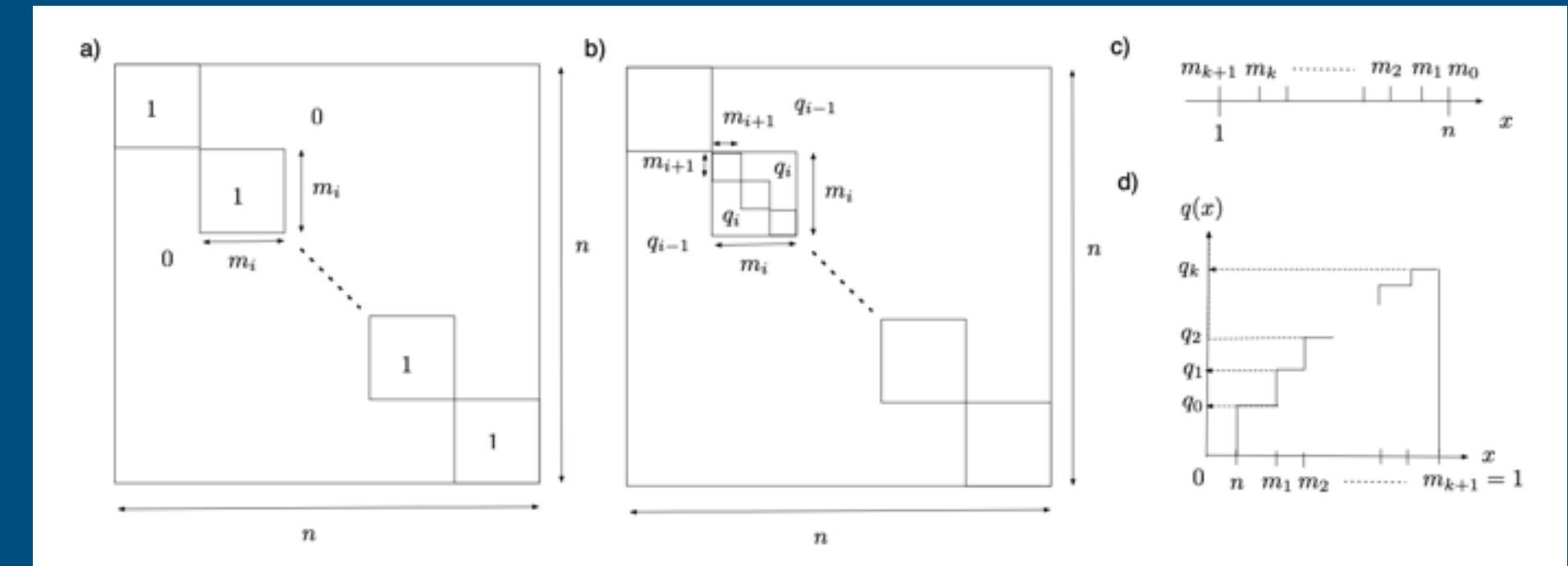
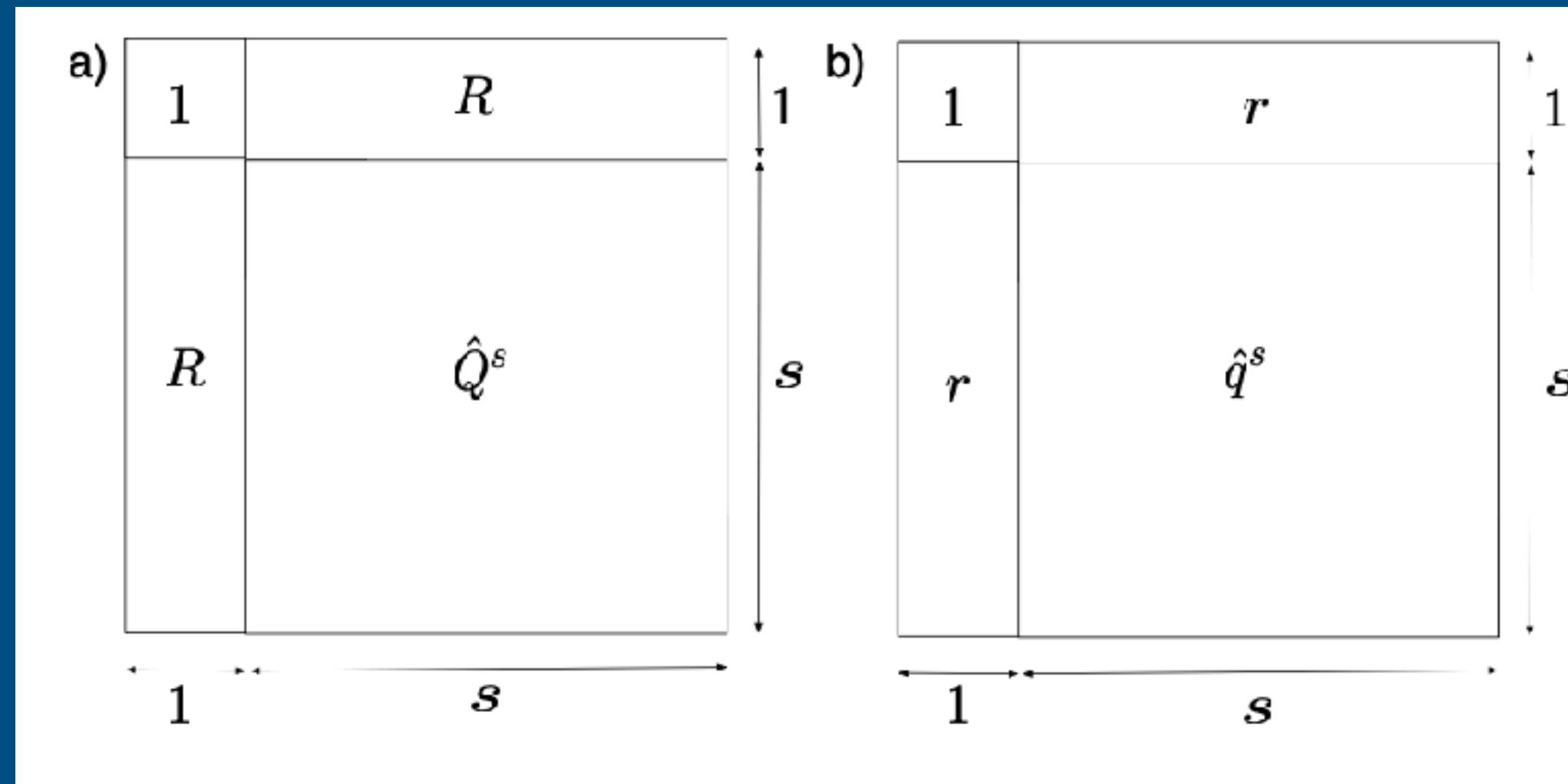
$$\frac{-\beta \overline{F(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} = \frac{\partial_n \overline{V^{1+n}(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}} \Big|_{n=0}}{NM} = \partial_n S_{1+n}[\{\hat{Q}(l), \hat{q}(l)\}] \Big|_{n=0}$$

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

$$\alpha = \frac{M}{N}$$

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

Parisi's RSB ansatz



$$Q_{ab}(l) = \sum_{k \neq 0}^{k+1} Q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L$$

$$q_{ab}(l) = \sum_{i=0} q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L - 1$$

Input/output boundaries

overlap among students
at the boundaries

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

overlap between the students
and the teacher

$$r(0) = r(1) = r$$

0

Bayes-optimal limit

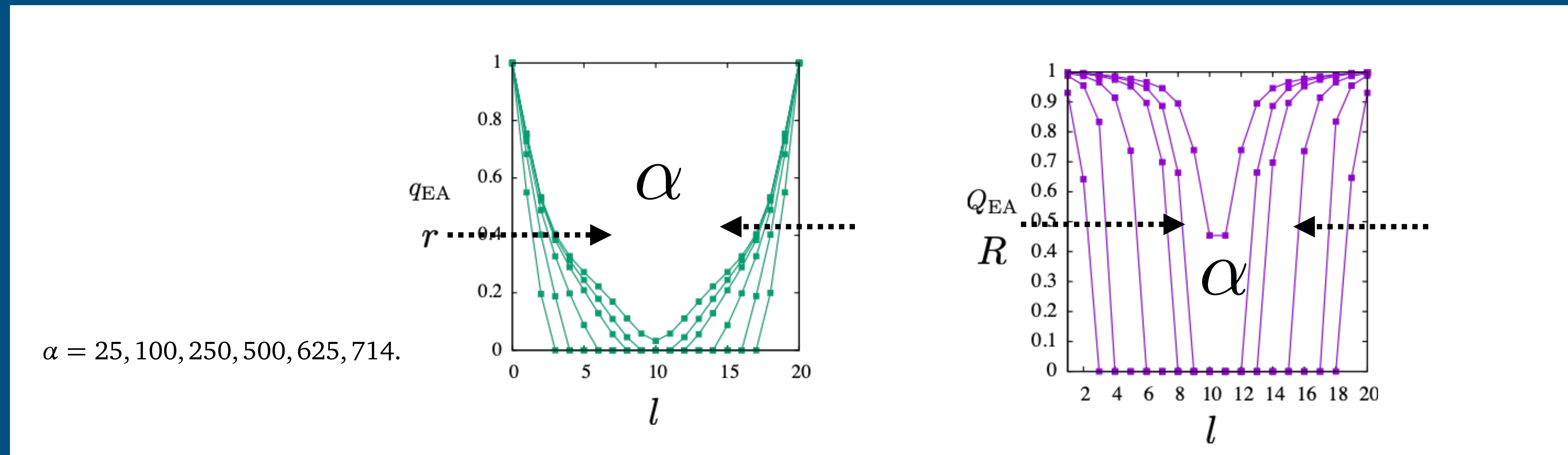


1

r

Bayes-optimal case (no noise)

Hajime Yoshino, SciPostPhys. Core 2, 005 (2020).



“Wetting transition”

layer-by-layer 2nd order transition

Bayes optimal, Nishimori condition $q = r, Q = R$

Replica symmetry (RS) holds

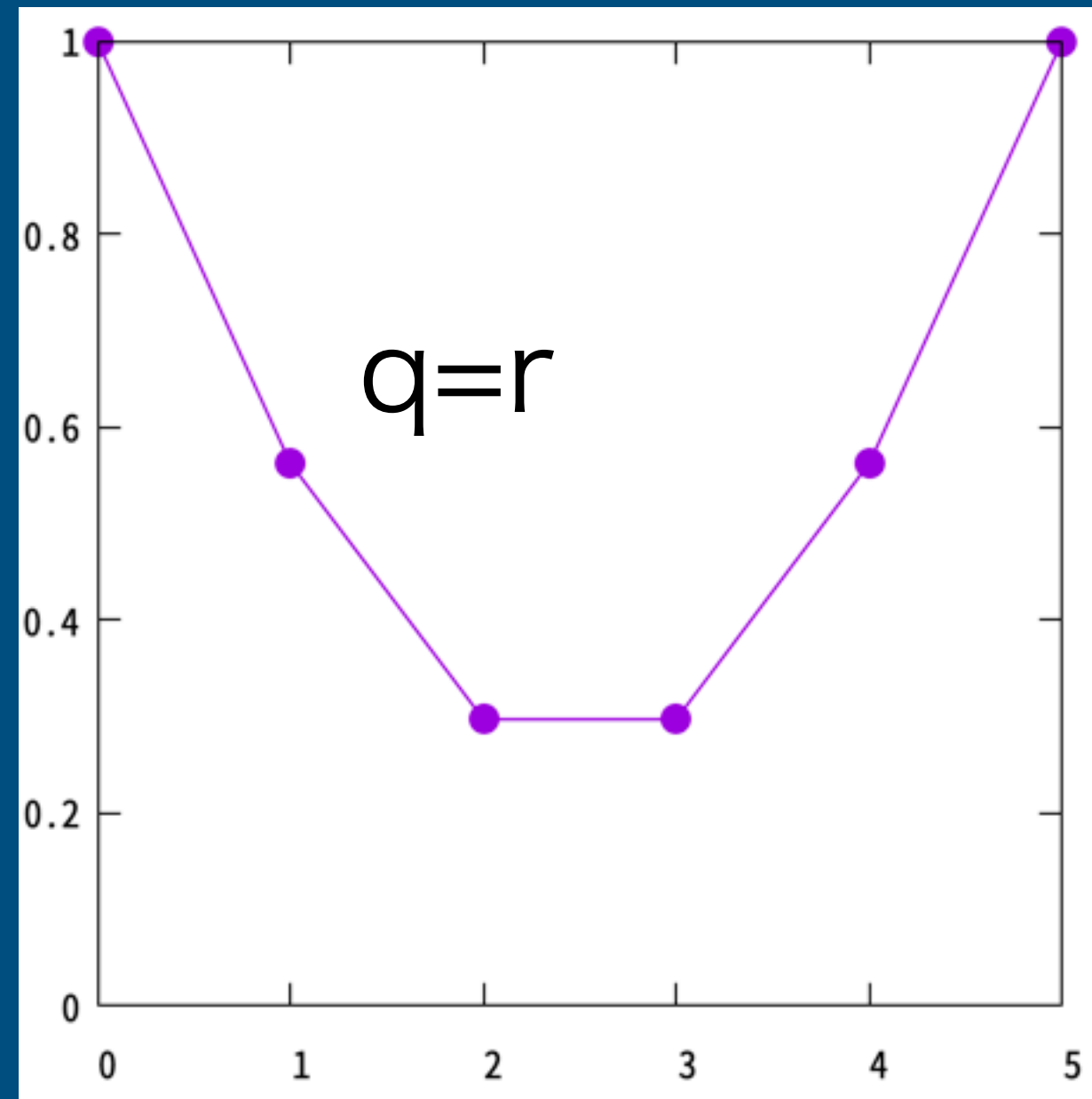
With noise

Spatial profile of the EA order parameter



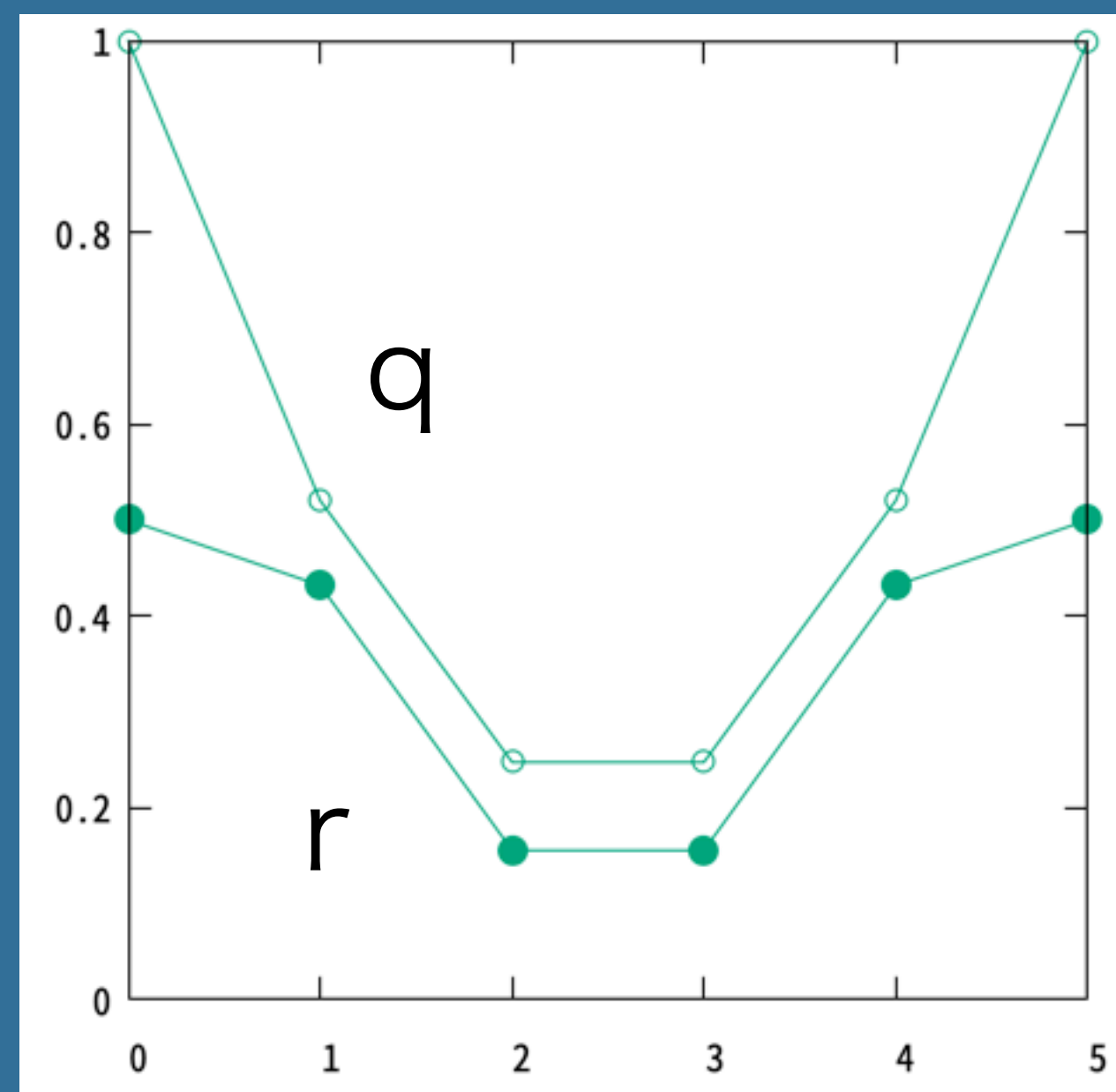
$$q_{EA} = q(1)$$

$$r = 1$$



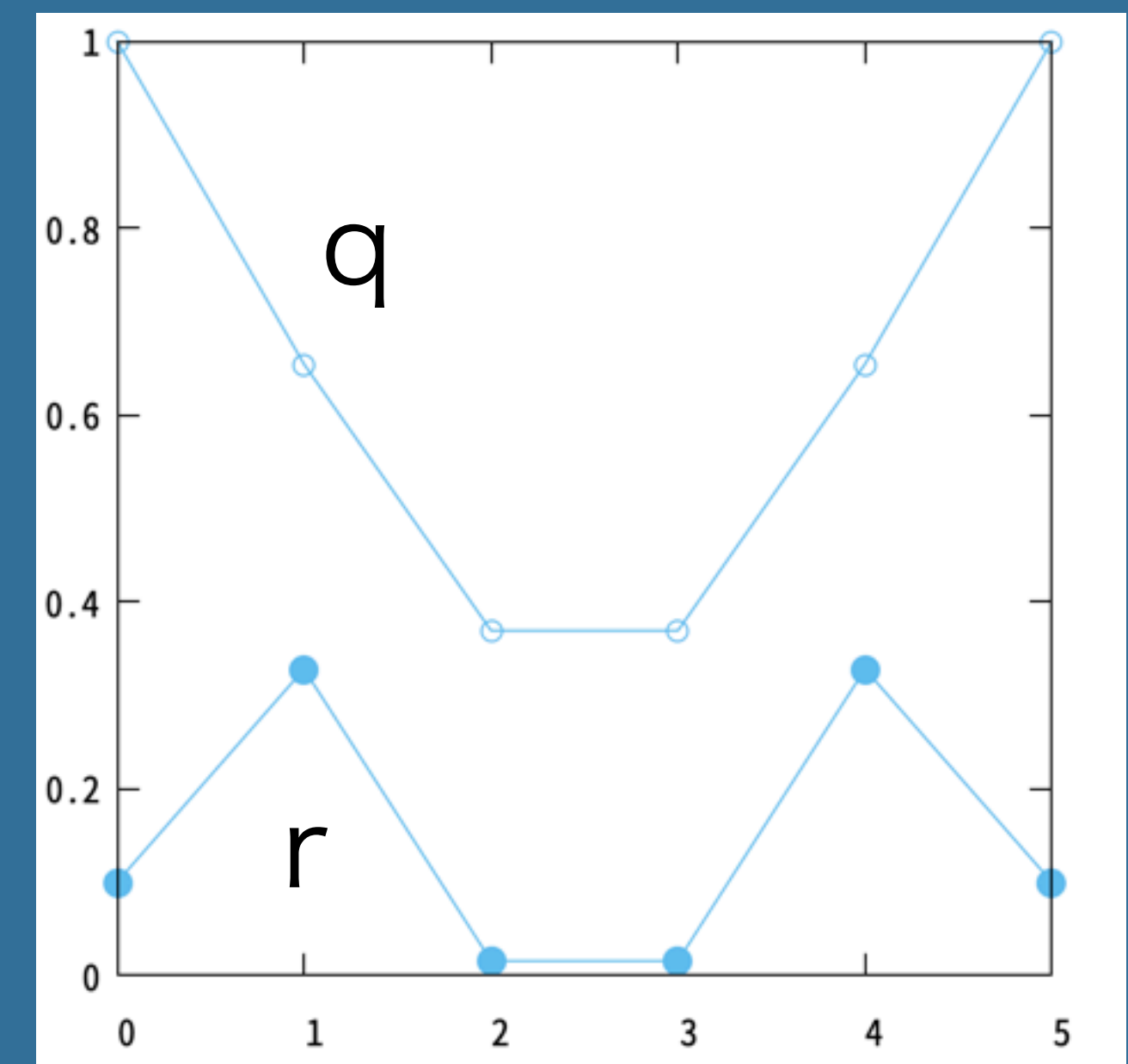
l

$$r = 0.5$$



l

$$r = 0.1$$



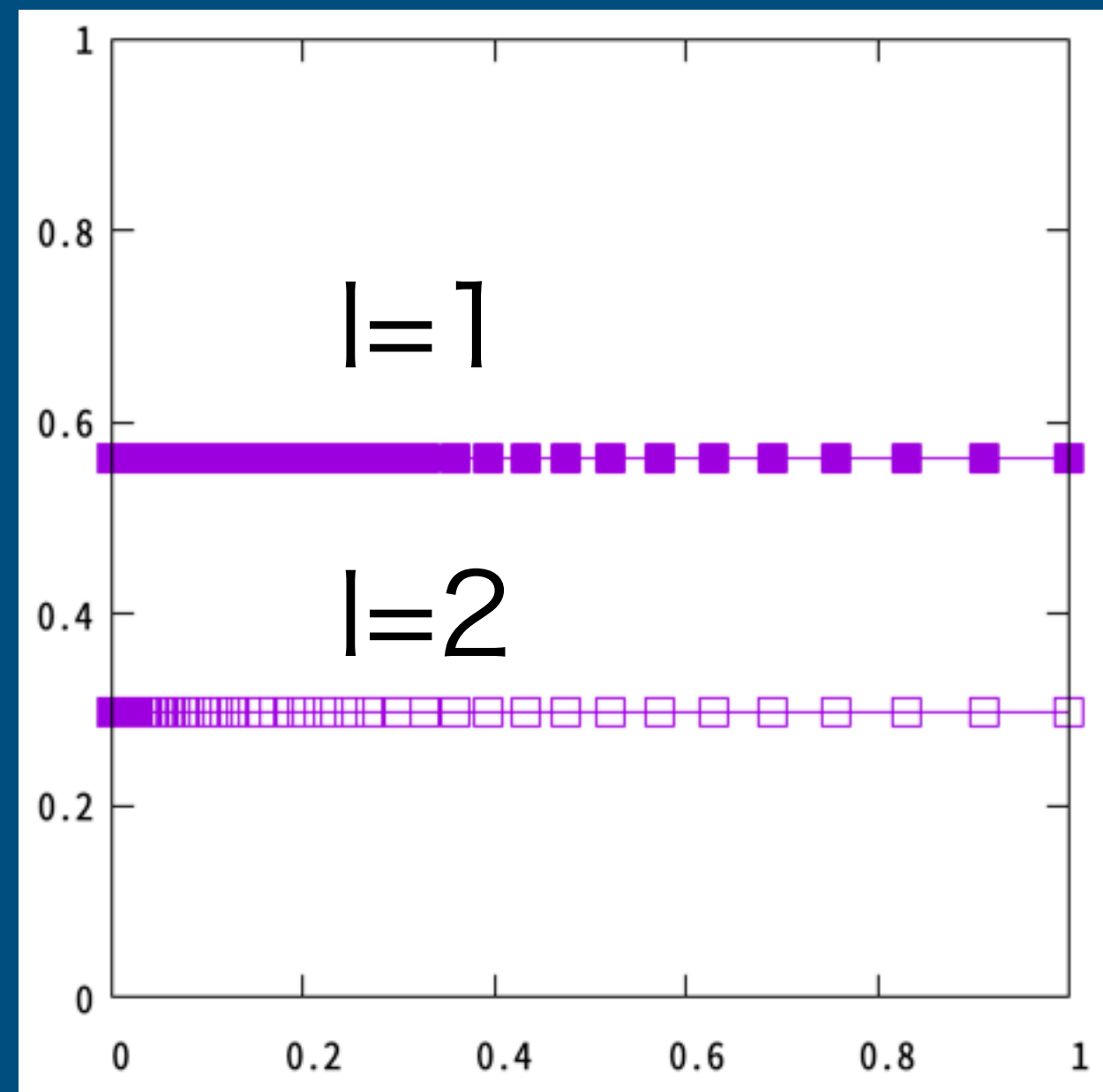
l

$$\alpha = 25$$

Hierarchical structure of the solutions

$$q(x)$$

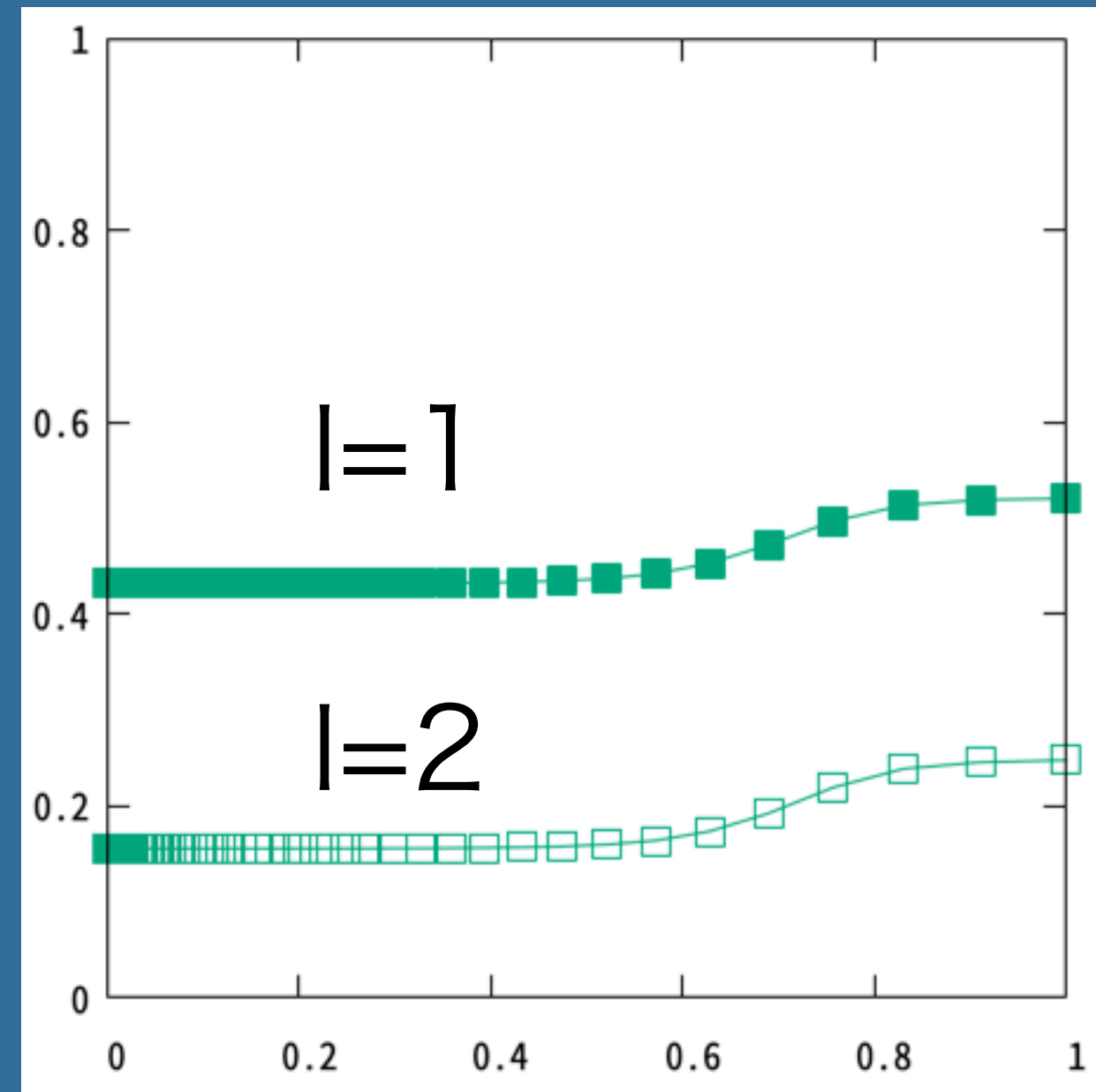
$$r = 1$$



x

Replica symmetric

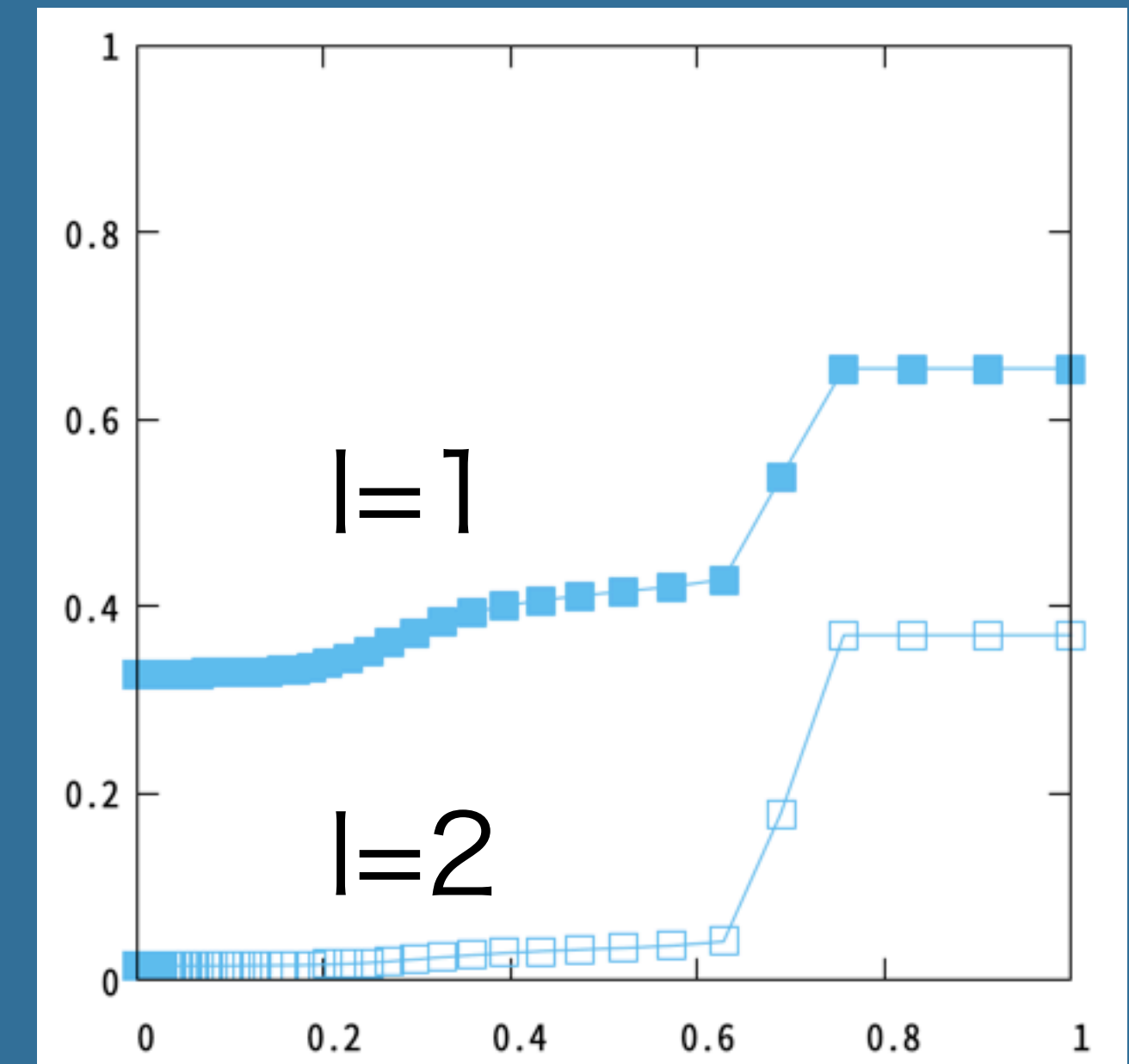
$$r = 0.5$$



x

Replica symmetry broken

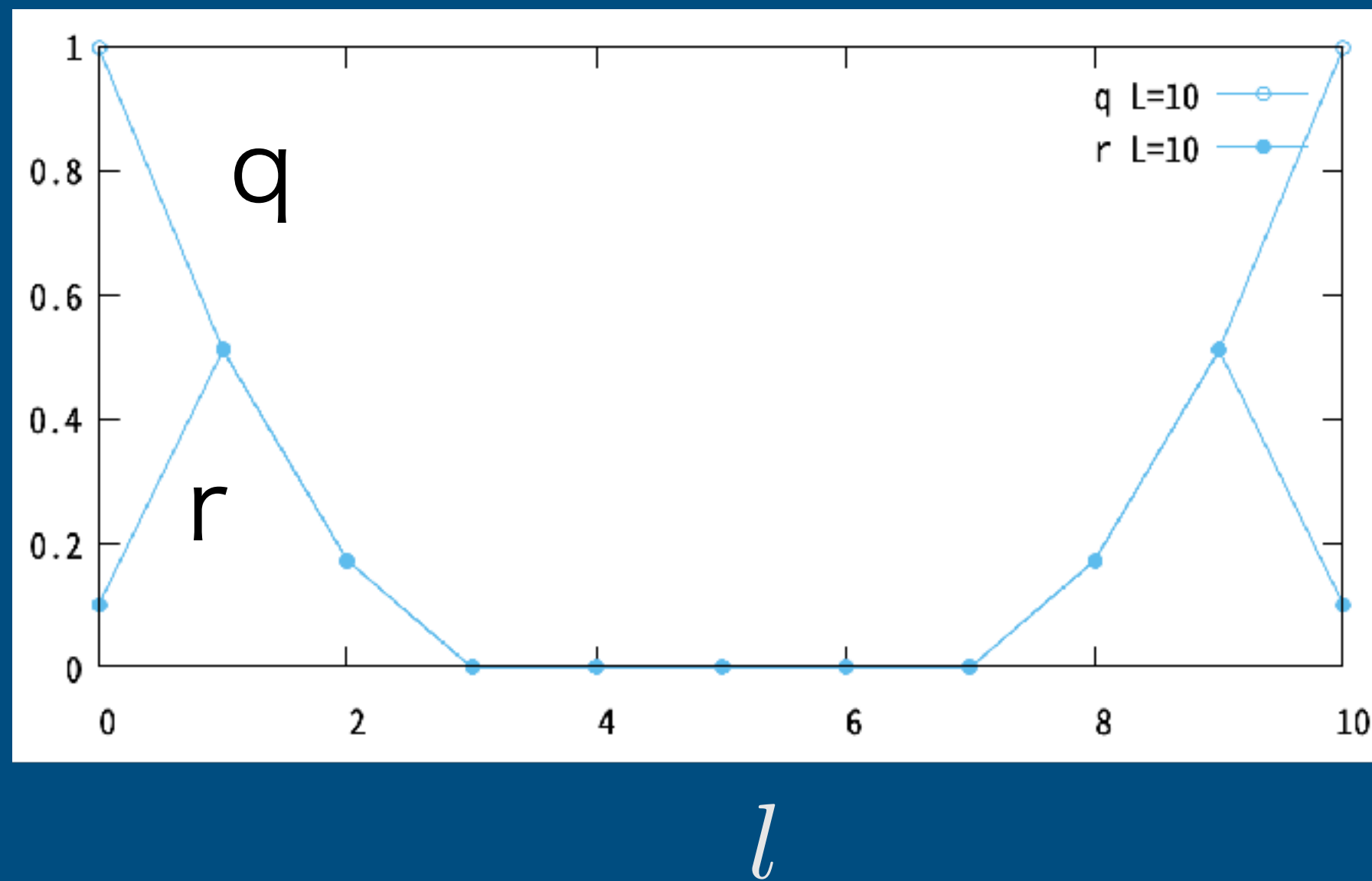
$$r = 0.1$$



x

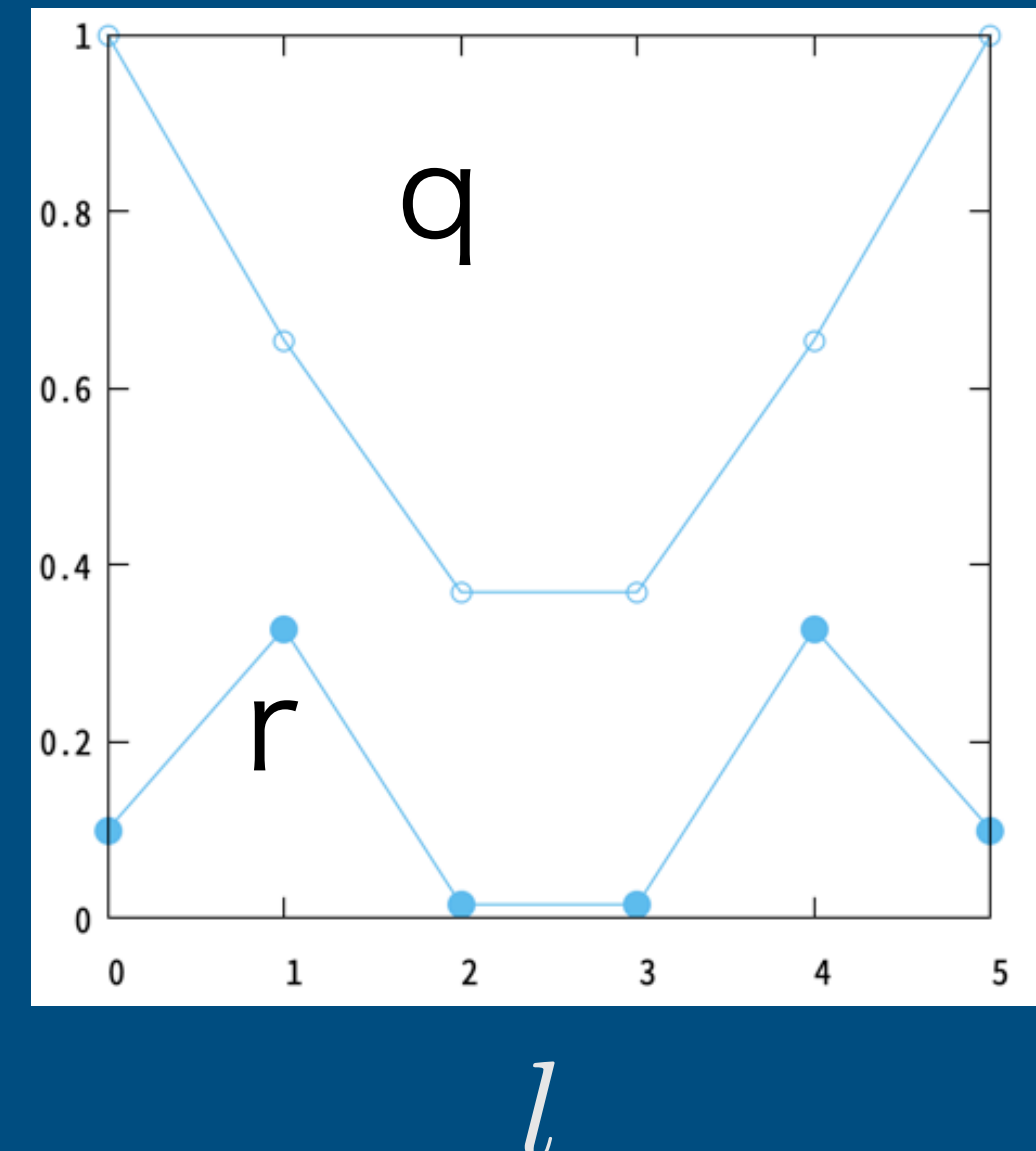
$r = 0.1$

$q_{EA} = q(1)$



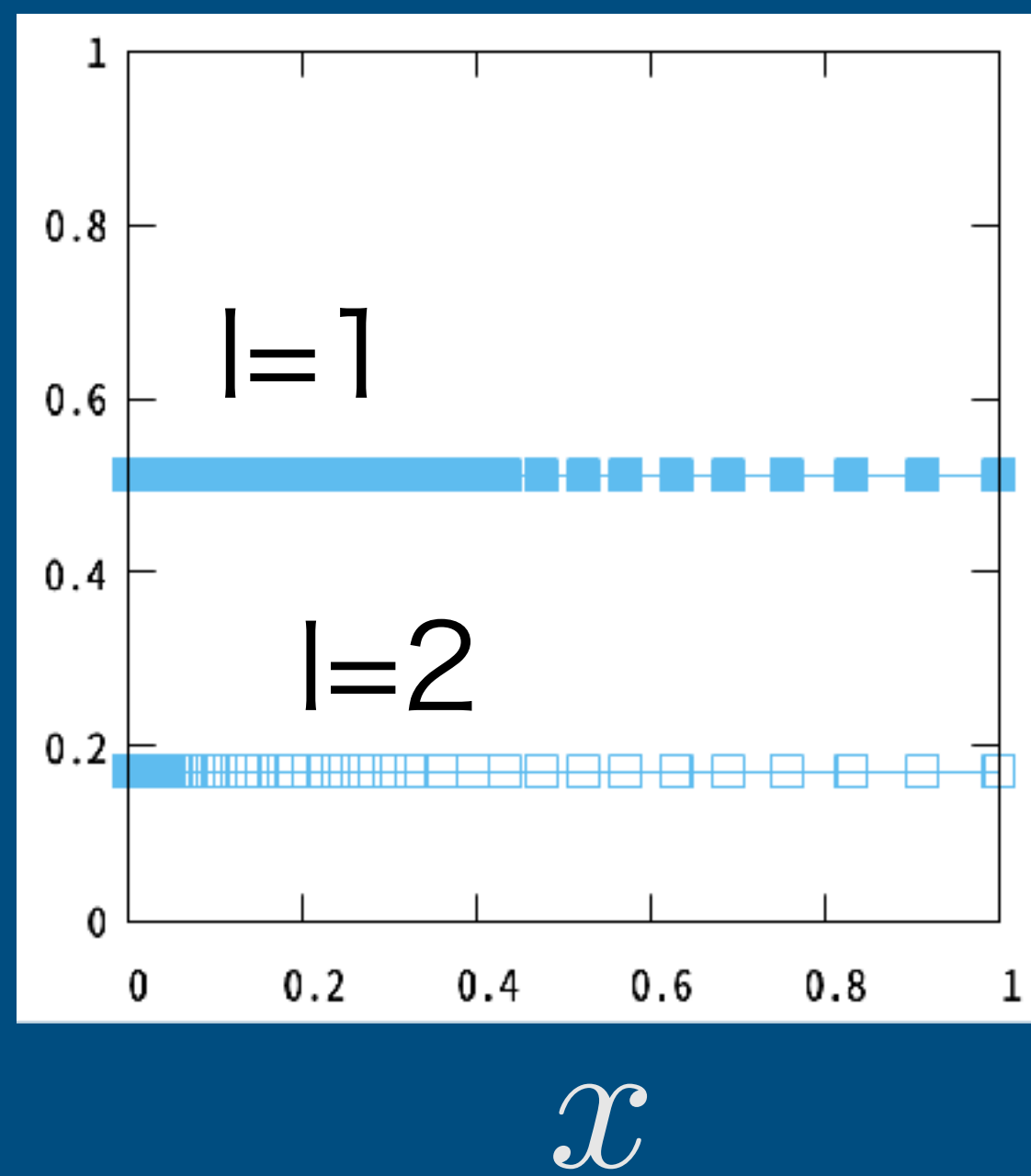
$r = 0.1$

$q_{EA} = q(1)$

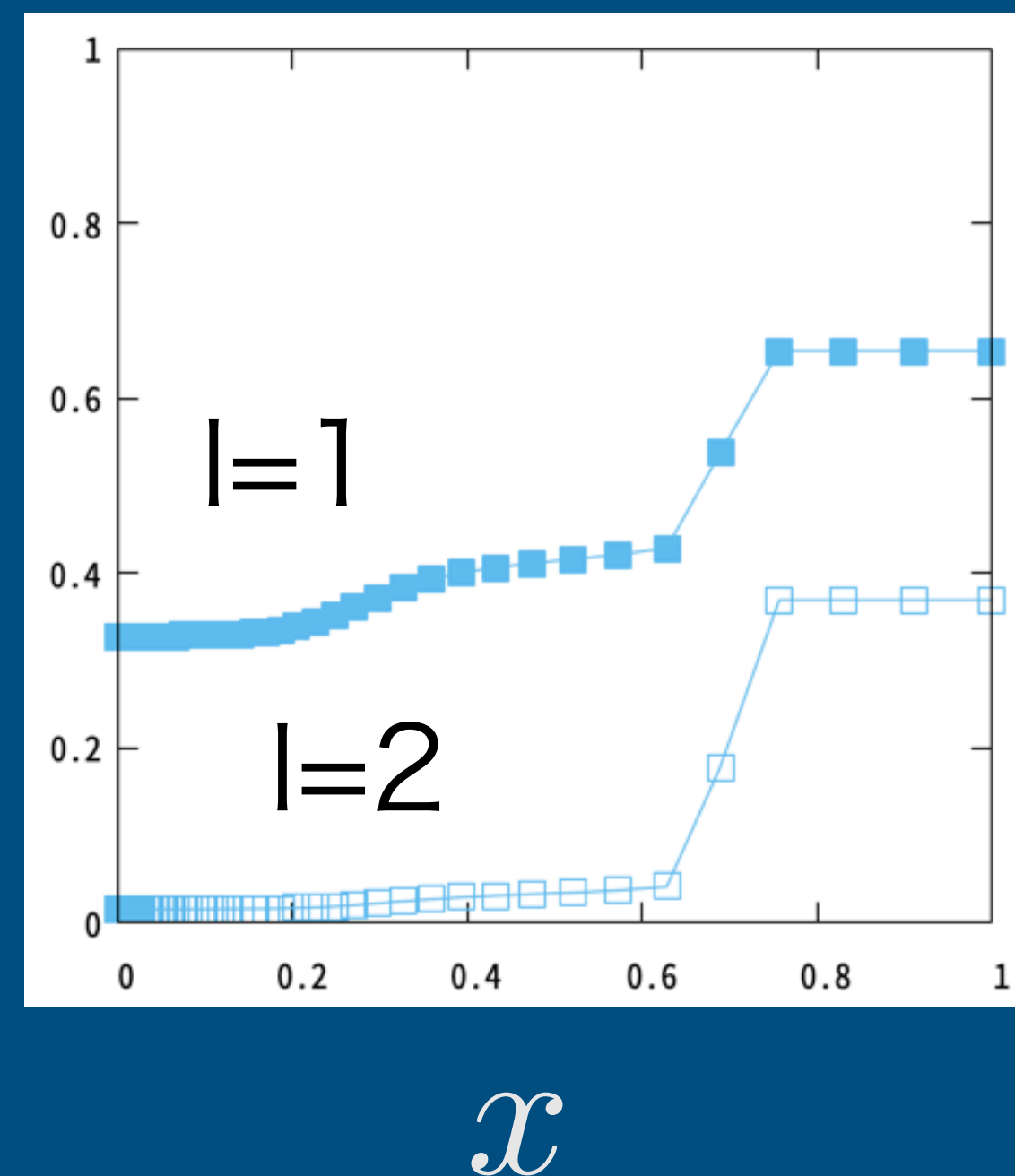


Replica symmetric!

$q(x)$



$q(x)$



■ Simulation of learning in a teacher-student setting

t

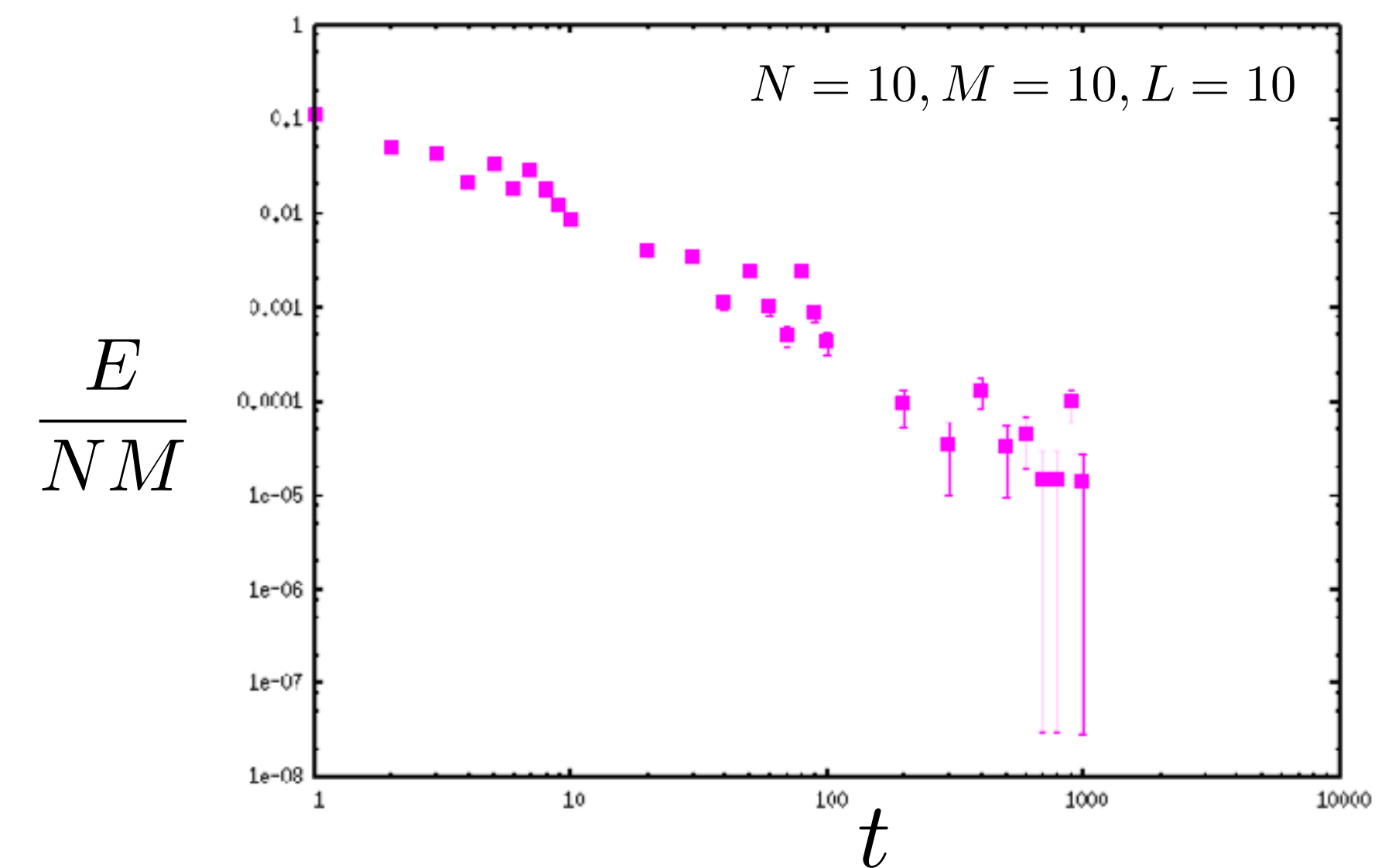
T=0 greedy Monte Carlo

random teacher + student “a” and “b”

loss function $E = \sum_{i=1}^N \sum_{\mu=1}^M \left(S_{L,i}^{\mu} - (S_*)_{L,i}^{\mu} \right)^2$

1. “**Unlearning**” : start from teacher’s configuration

2. “**Learning**” : start from random configuration

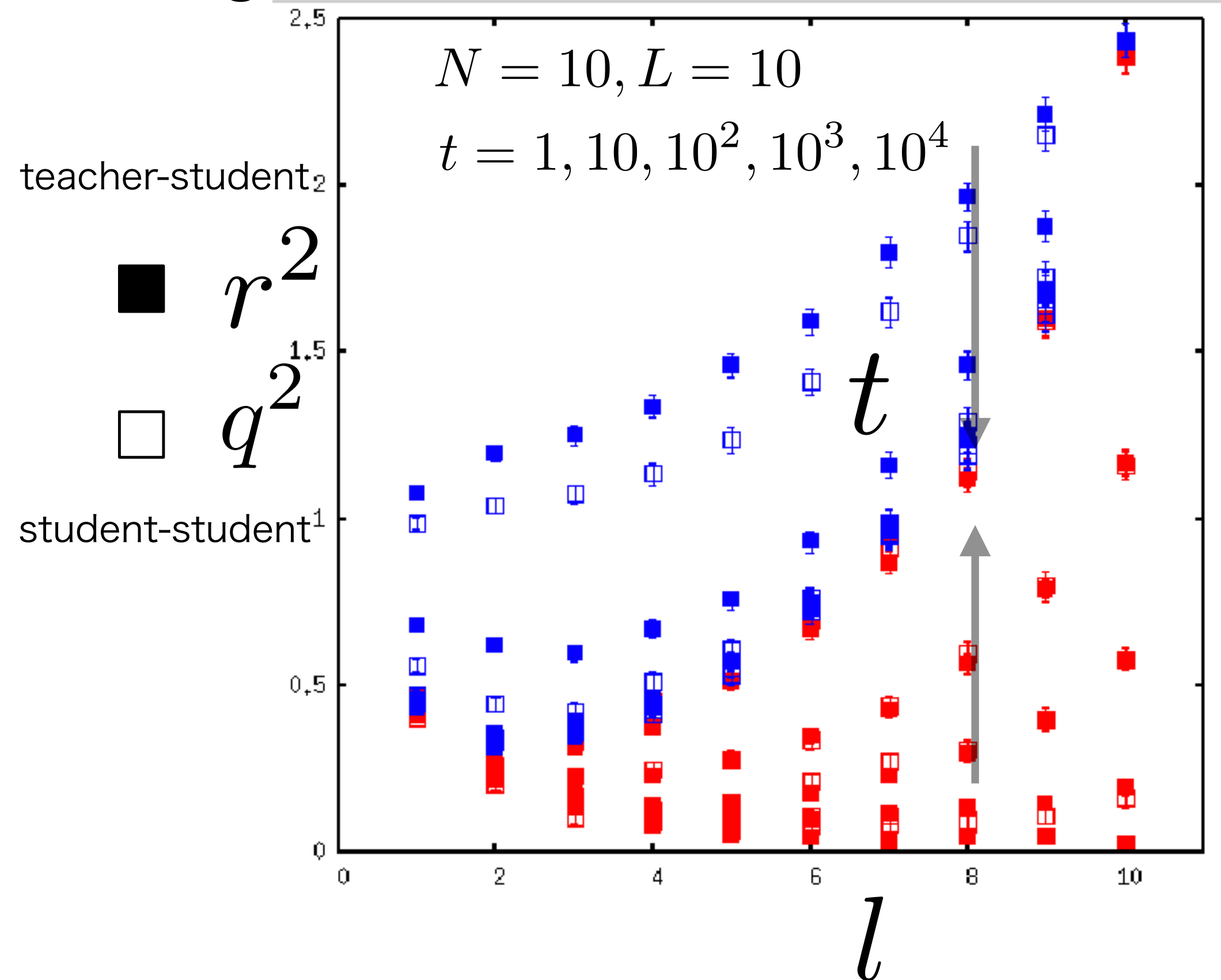


Permutation-invariant overlap

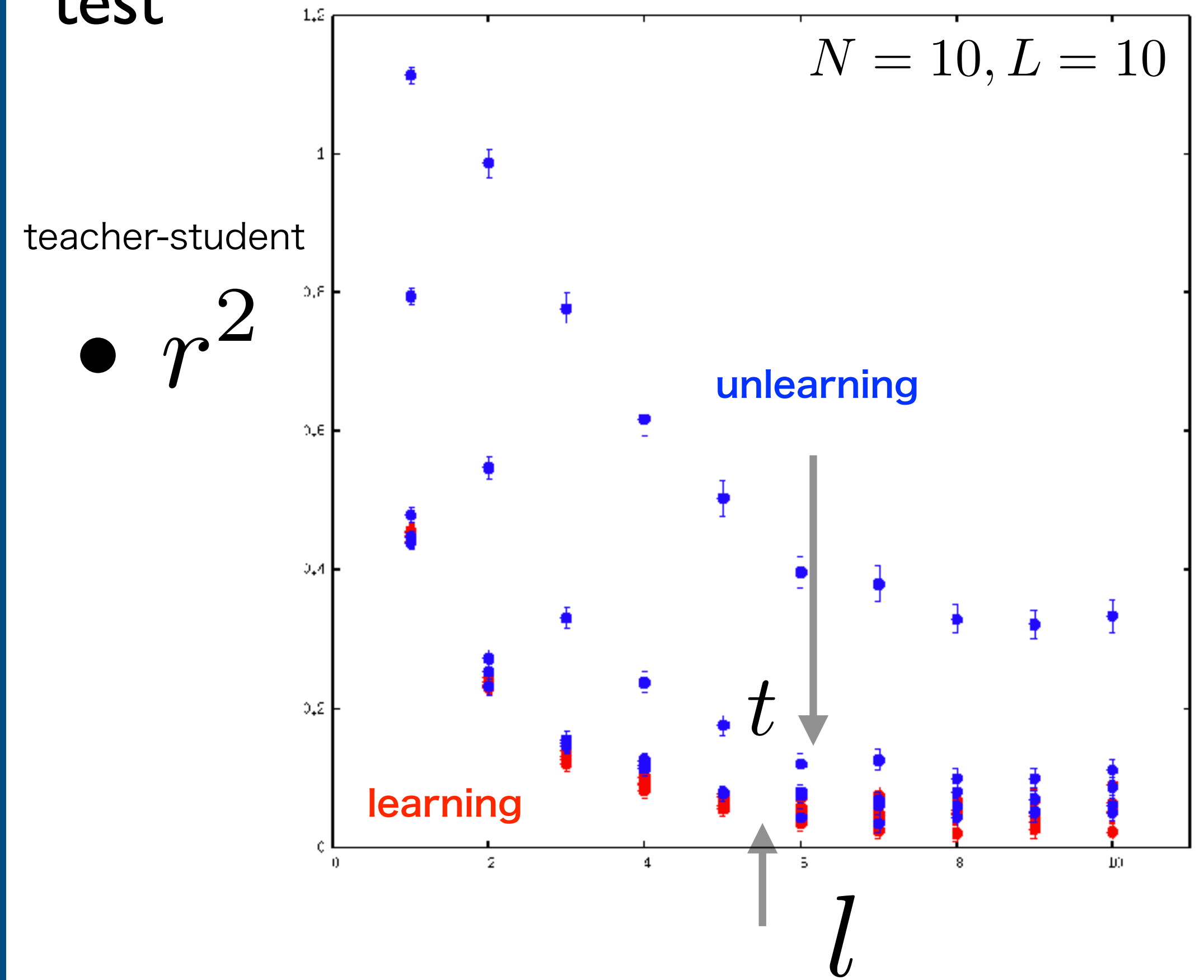
$$q^2 \equiv \frac{1}{N} \sum_{i,j=1}^N q_{ij}^2 - \frac{1}{\alpha} \quad q_{ij} = \frac{1}{M} \sum_{\mu=1}^M (S_i^\mu)^a (S_j^\mu)^b$$

$$\alpha = 1$$

training



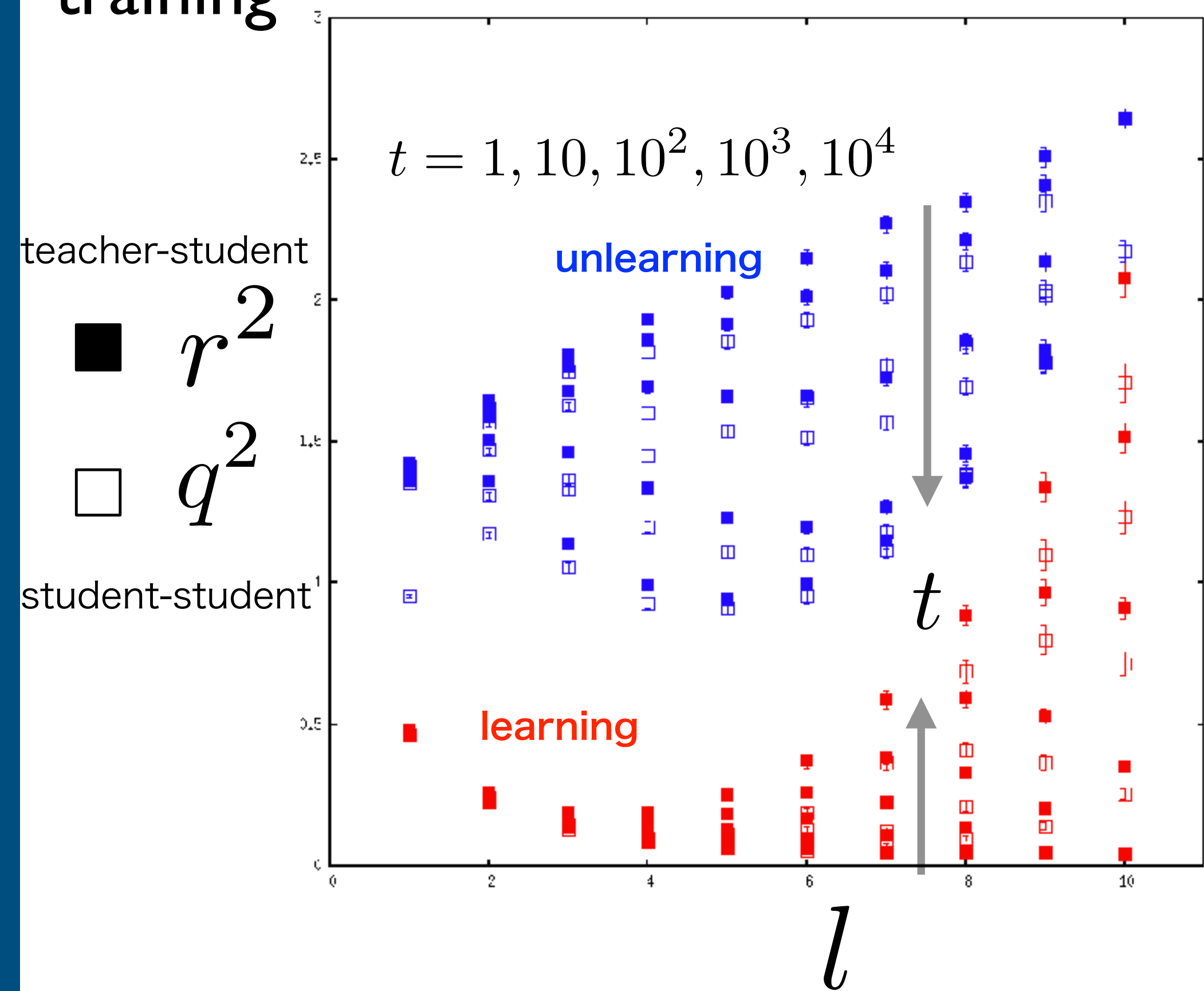
test



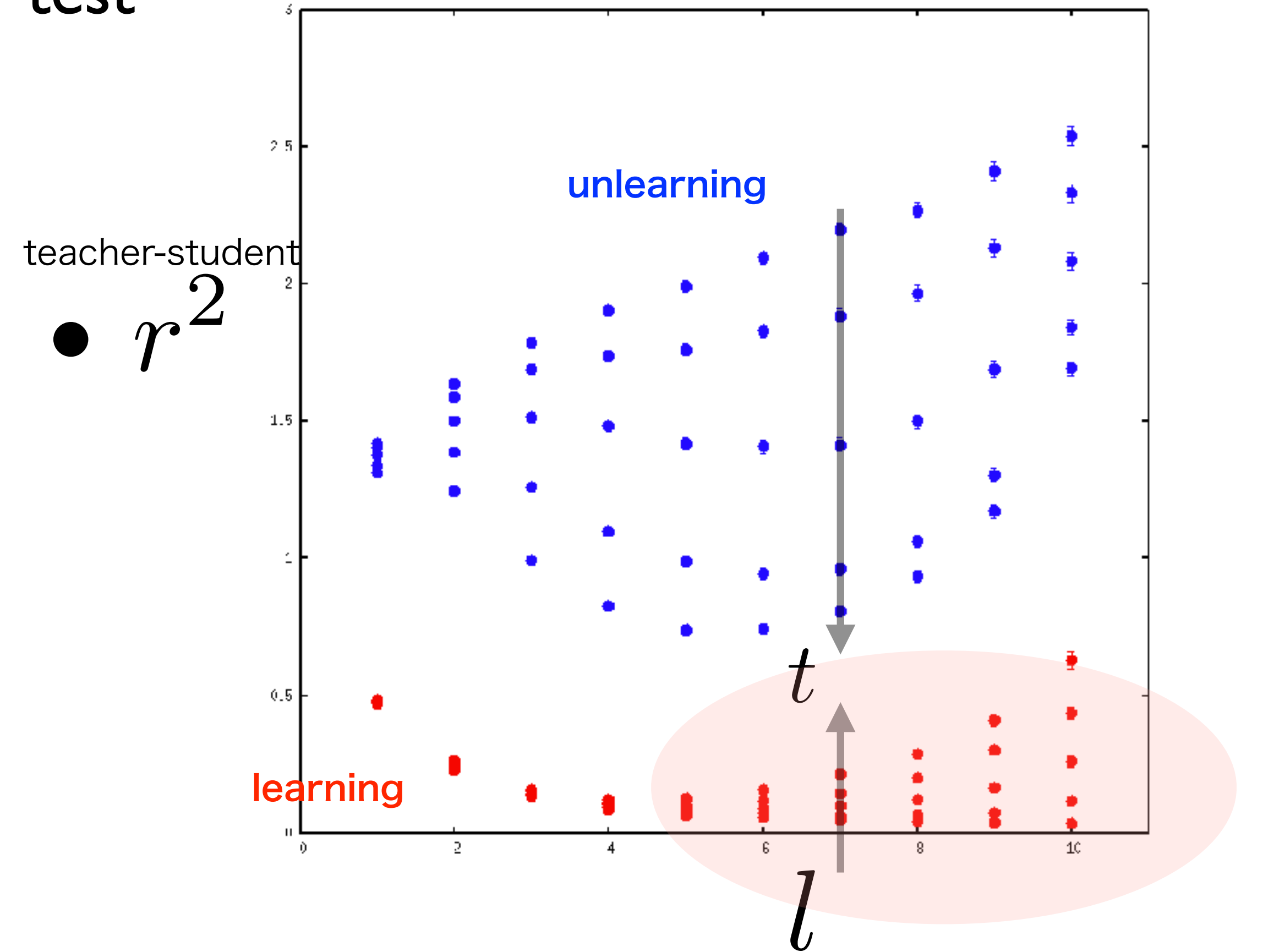
$$\alpha = 32$$

 $N = 10, L = 10$

training

 $N = 10, L = 10$

test

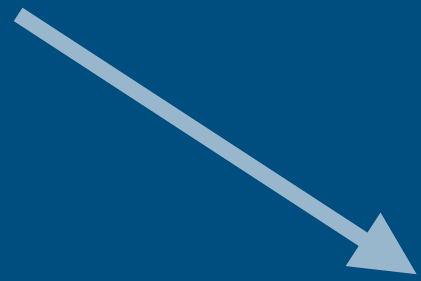


Generalization ability

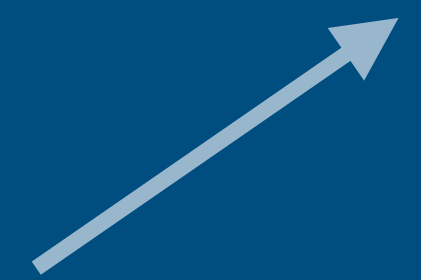
teacher-student overlap in the test

$$r_{\text{test}} = \frac{1}{NM} \sum_{i=1}^N \sum_{\mu=1}^M (S_i^\mu)^{\text{teacher}} (S_i^\mu)^{\text{student}}$$

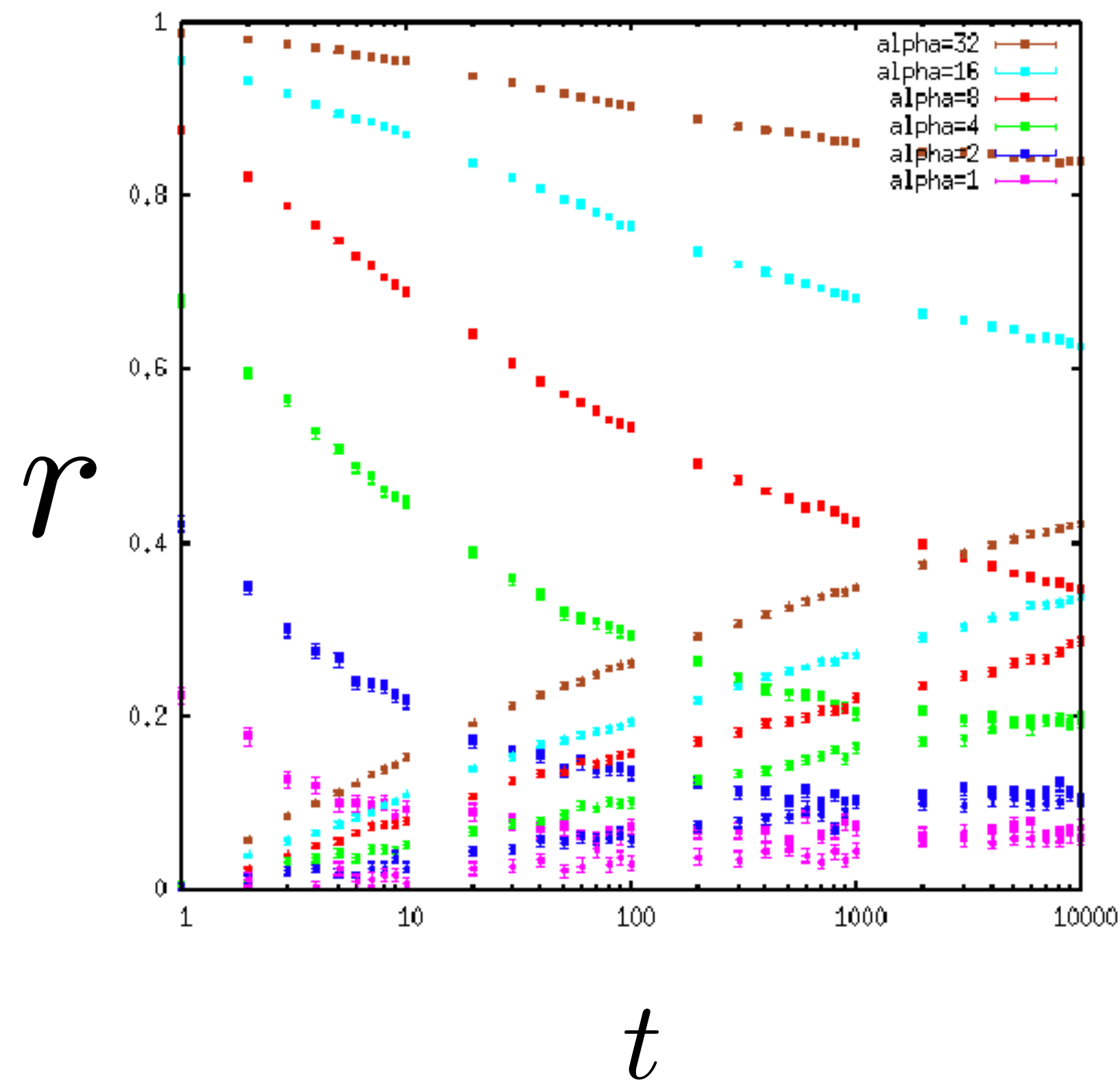
unlearning



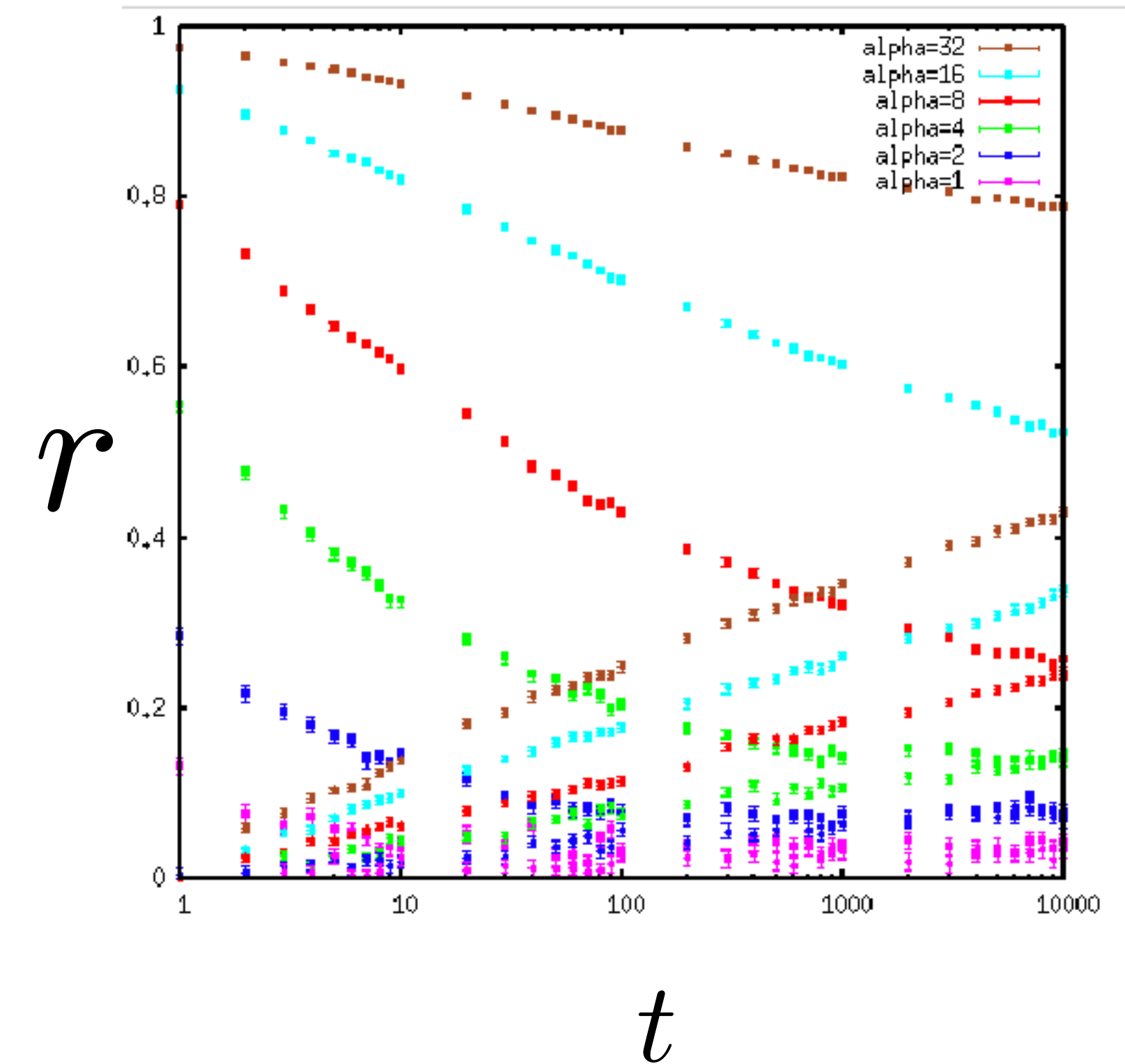
learning



$N = 10, L = 5$

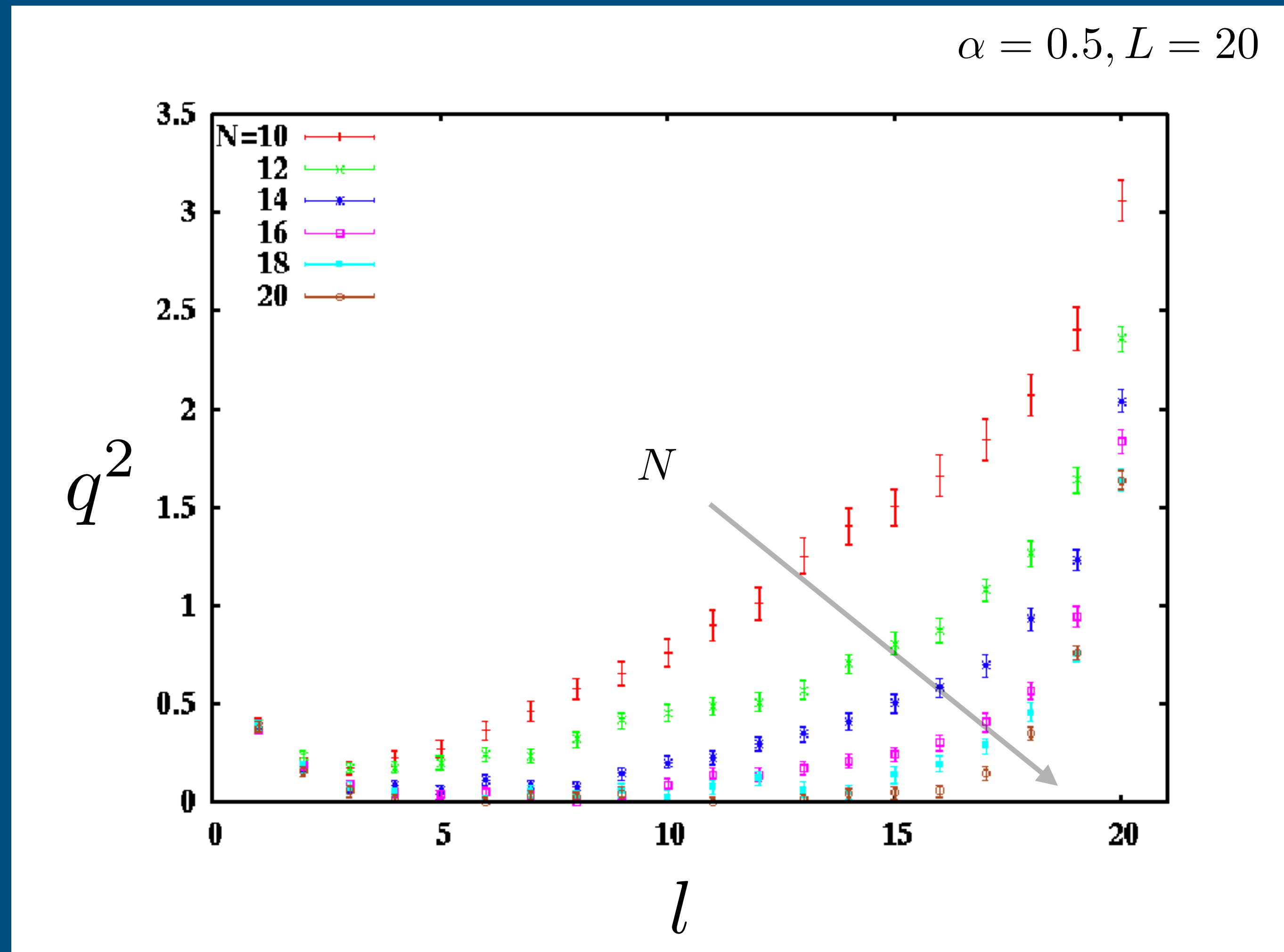


$N = 10, L = 10$



deeper systems generalize as well...

Increasing N with fixed $\alpha = M/N$



more symmetric with larger N

■ Summary

Construction of **replica theory** for a deep perceptron network

- random input/output (random constraint satisfaction problem)
- teacher-student scenario (statistical inference) **with noise**

“**Wetting transition**” in the design space with/without **RSB**

Numerical simulations of the teacher-student scenario

■ Outlook

Goldt, S., Mézard, M., Krzakala, F., & Zdeborová, L. (2020). PRX, 10(4), 041044.

Finite width N effect, hidden manifold model

c.f. MNIST $N = 784$ but $D_{\text{eff}} \simeq 14$

mismatch of architecture

other activation functions: sigmoid, ReLU, ...

Simulations with “real data”, various algorithms, architectures...