

深層ニューラルネットワークにおける レプリカ対称性の破れとその空間構造

Hajime Yoshino

Cybermedia Center,Osaka University



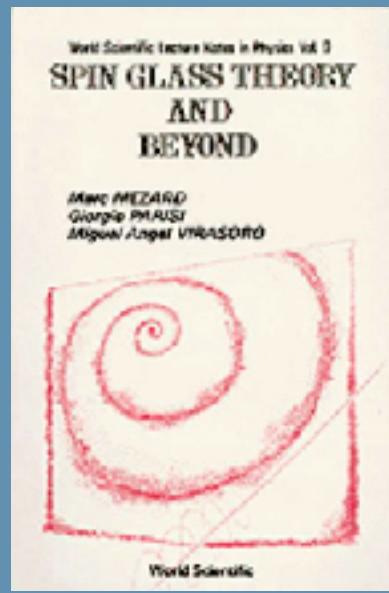
Hajime Yoshino, SciPostPhys. Core 2, 005 (2020).

「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

statistical mechanics of disordered systems



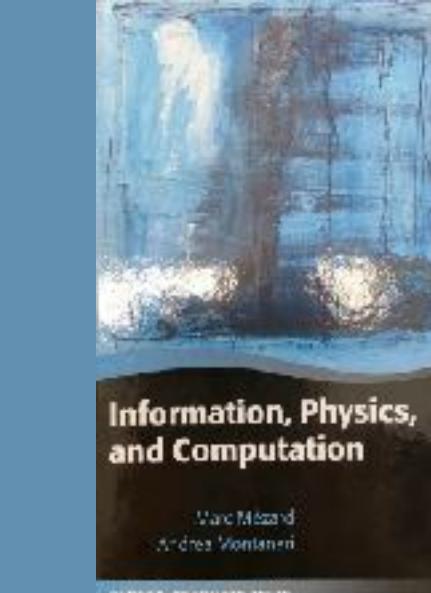
1987



1991

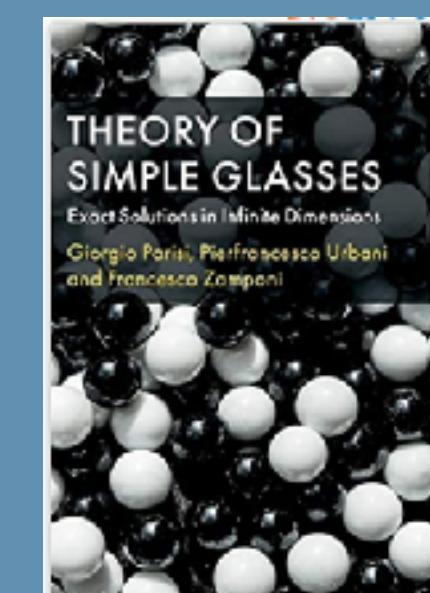


1999

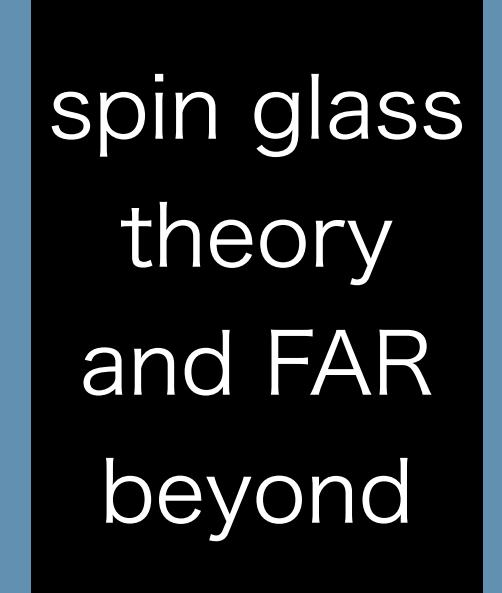


2001

without quenched disorder



2020



in progress

Statistical mechanics of disordered systems: spins, spheres and machines

Hajime Yoshino^{1,2}

¹Cybernetia Center, Osaka University, Toyonaka, Osaka 560-0043, Japan

²Graduate School of Science, Osaka University, Toyonaka, Osaka 560-0043, 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.

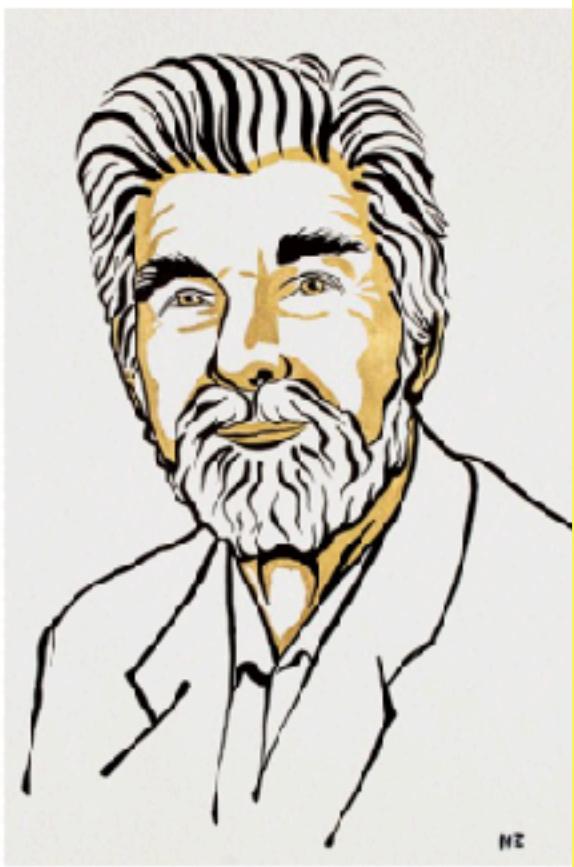
The Nobel Prize in Physics 2021



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



Disordered serendipity: a glassy path to discovery

A workshop in honour of Giorgio Parisi's 70th birthday
Sapienza Università di Roma, September 19-22, 2018

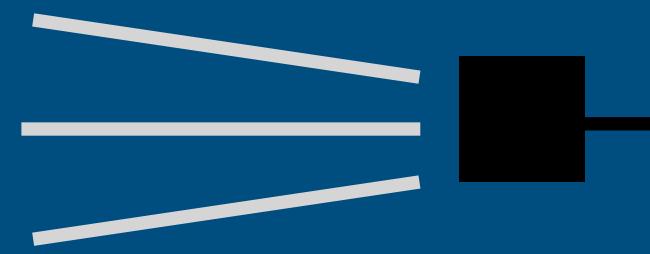


"for groundbreaking contributions to our understanding of complex physical systems"

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

Perceptron

Example

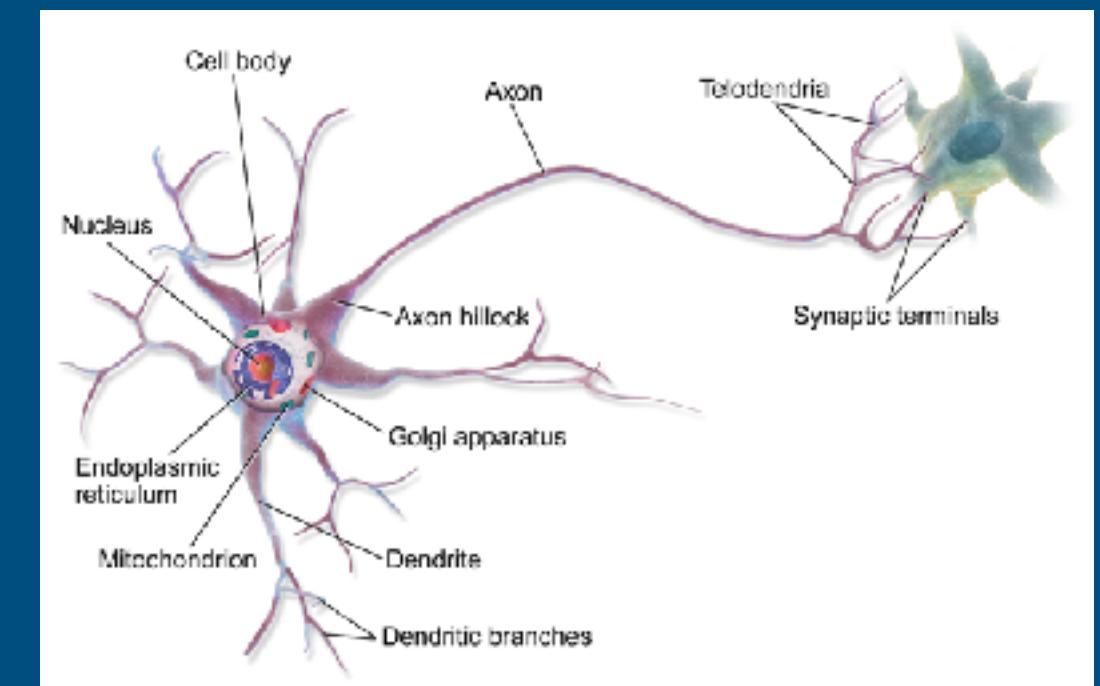
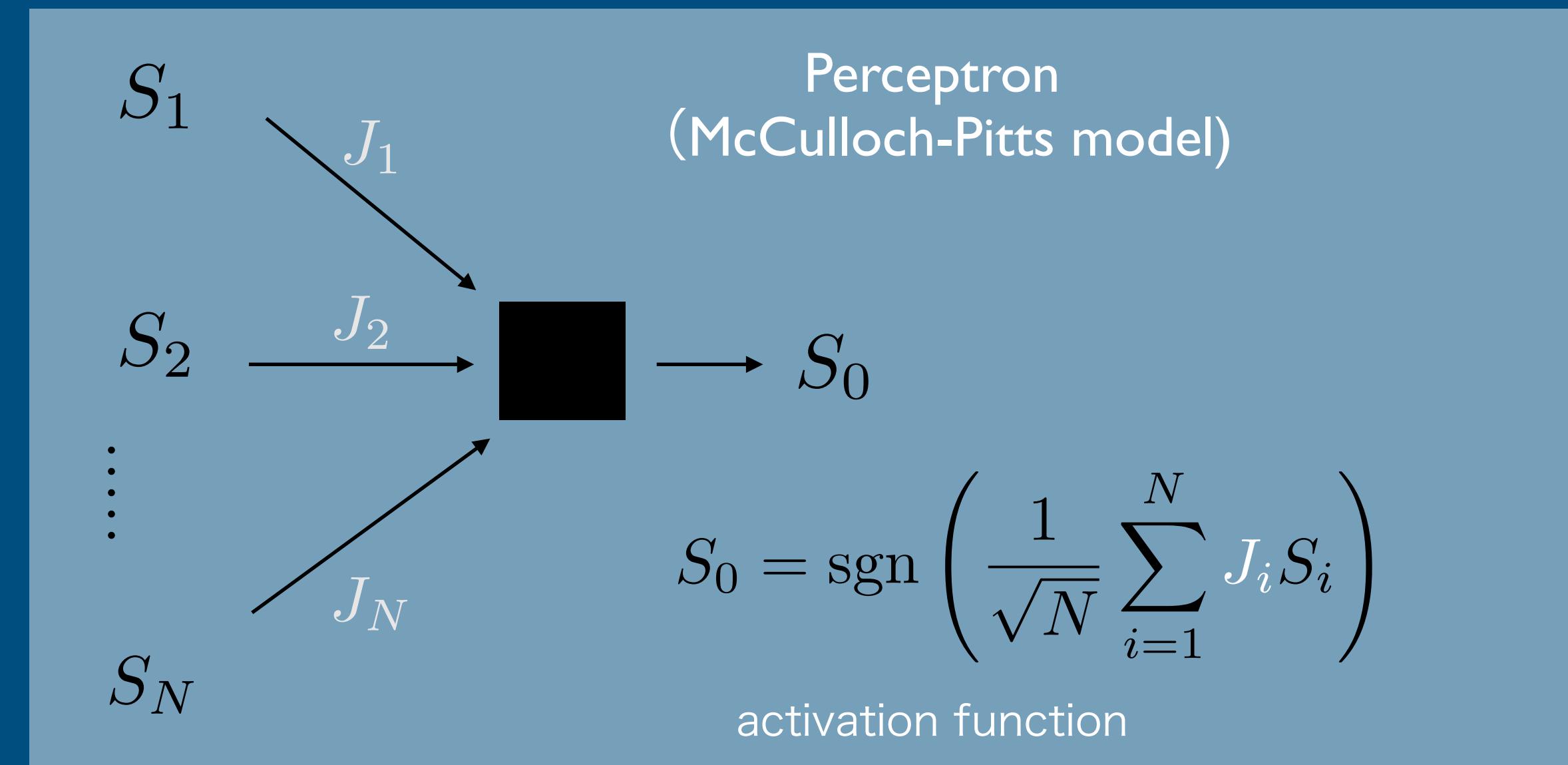


$$\begin{aligned} \text{\# of possible machines} \\ 2^3 = 8 \end{aligned}$$

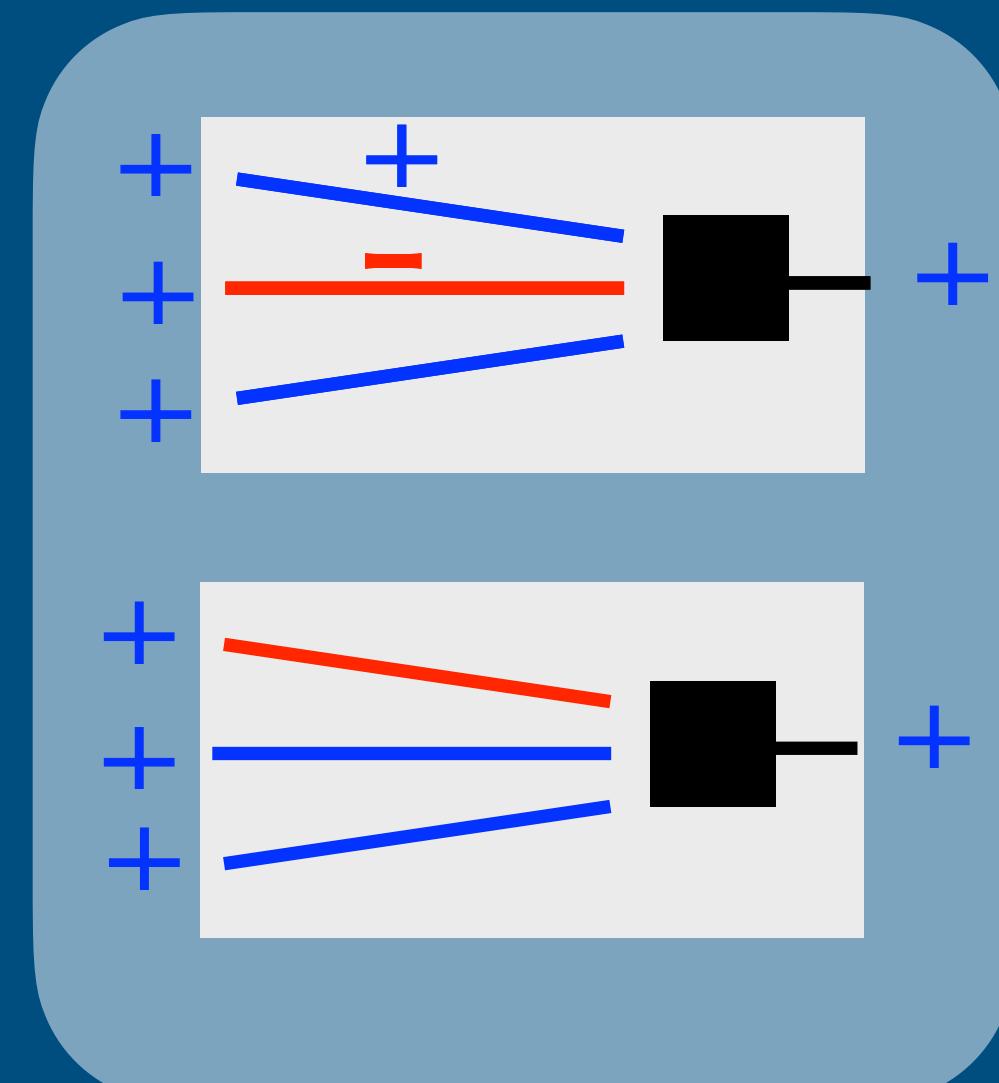
machine1

machine2

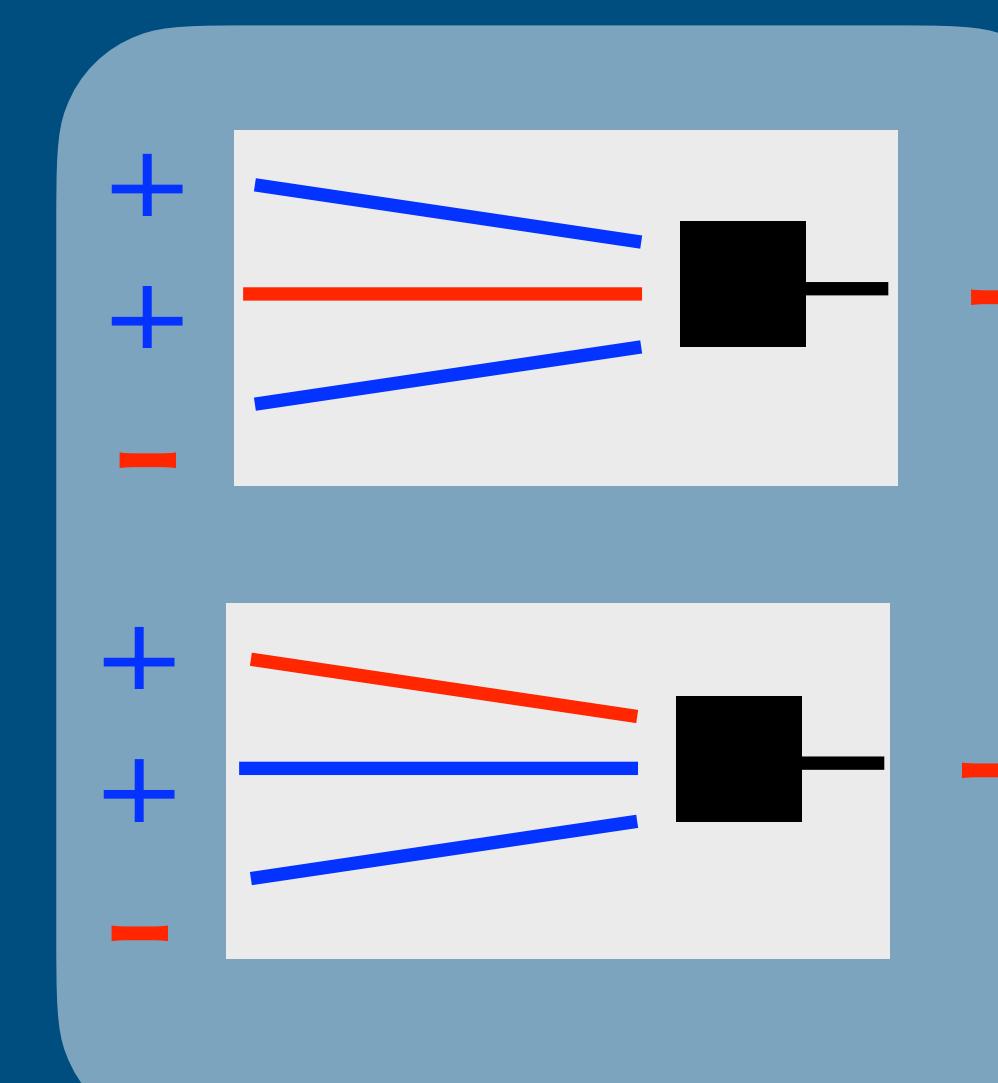
⋮



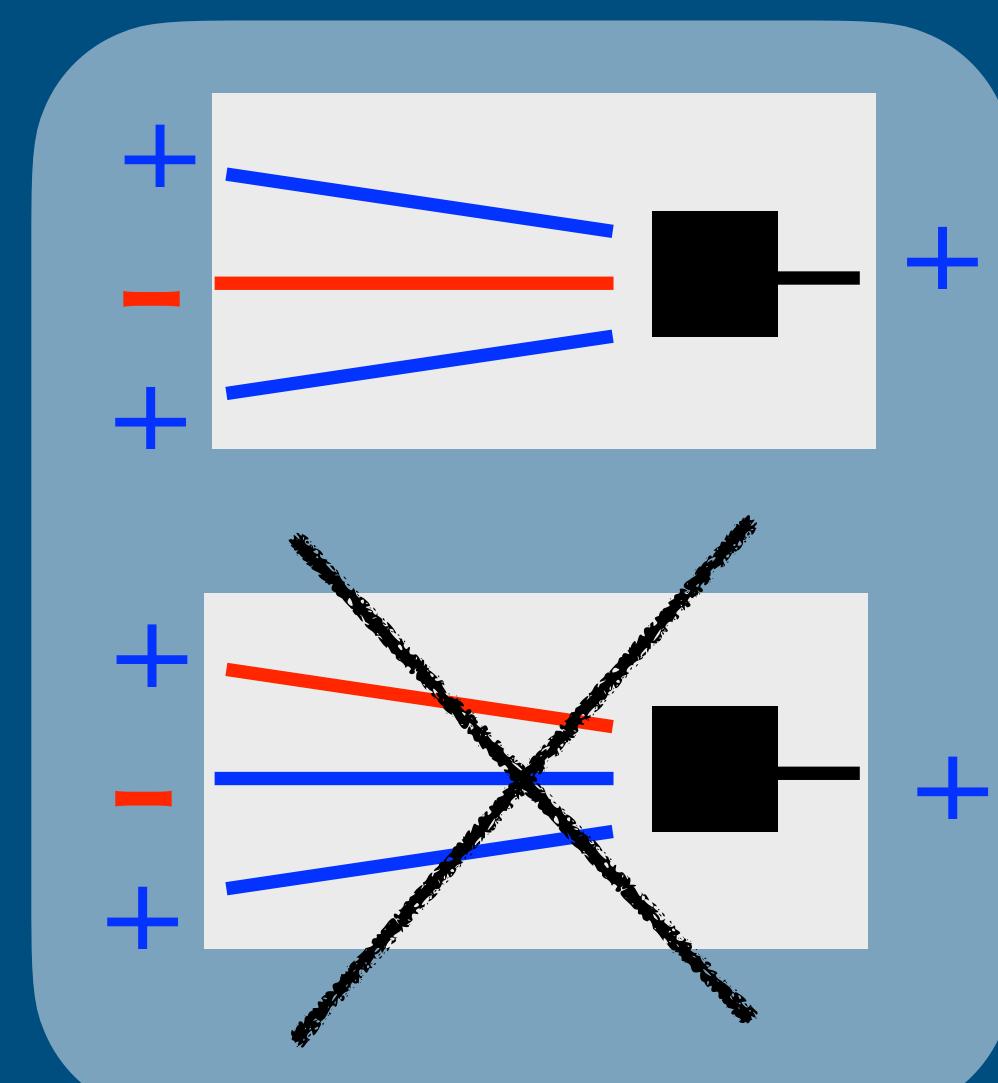
DATA1



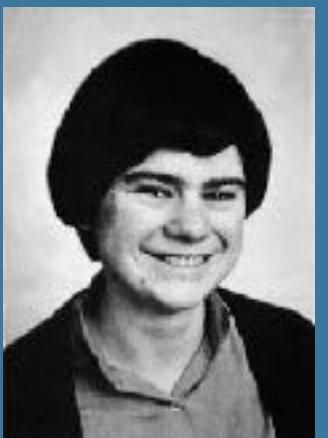
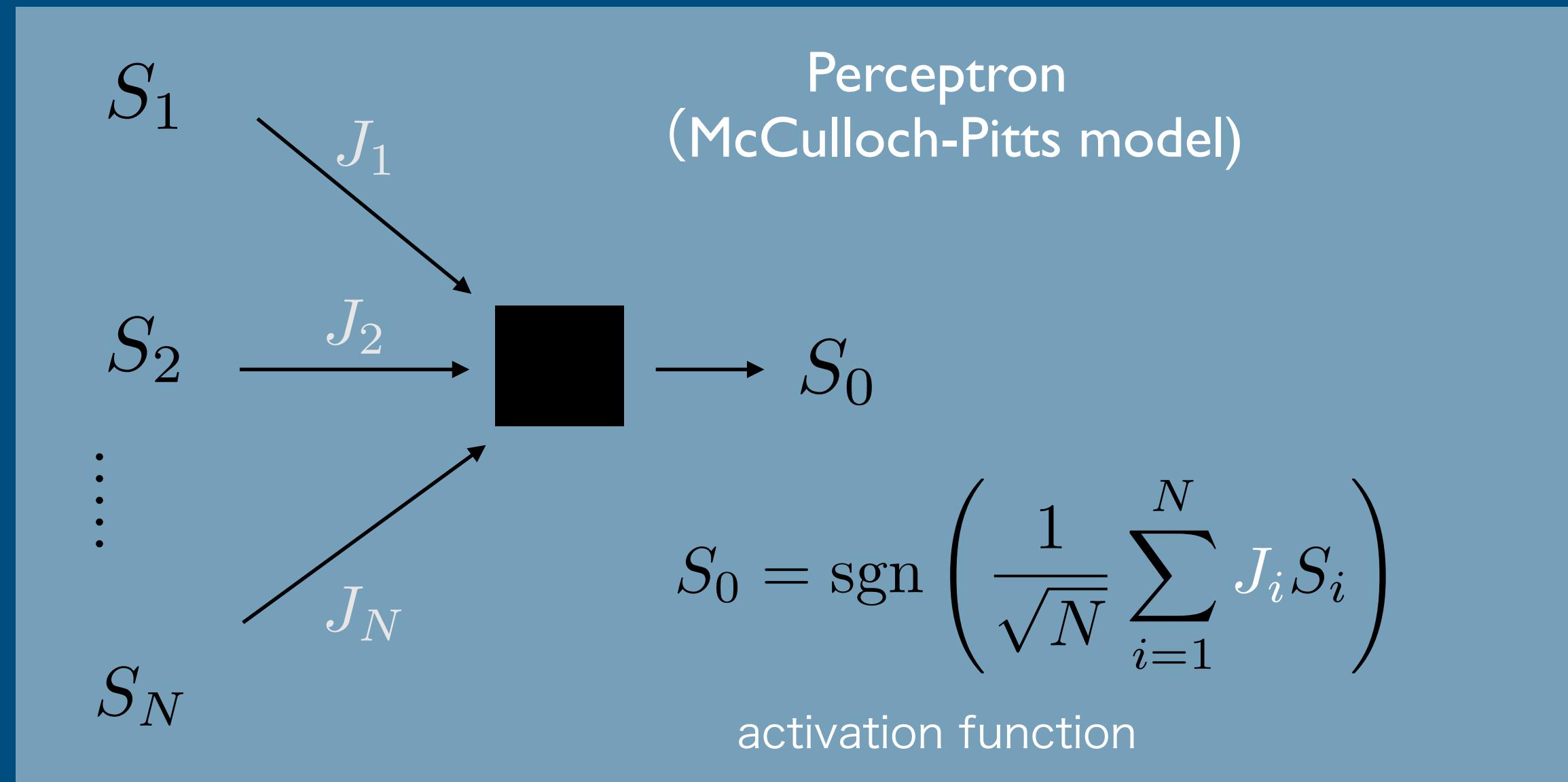
DATA2



DATA3



.....

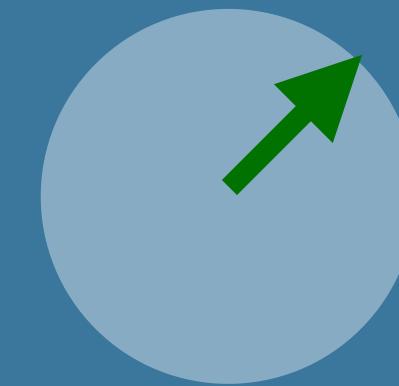


Elisabeth Gardner

(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

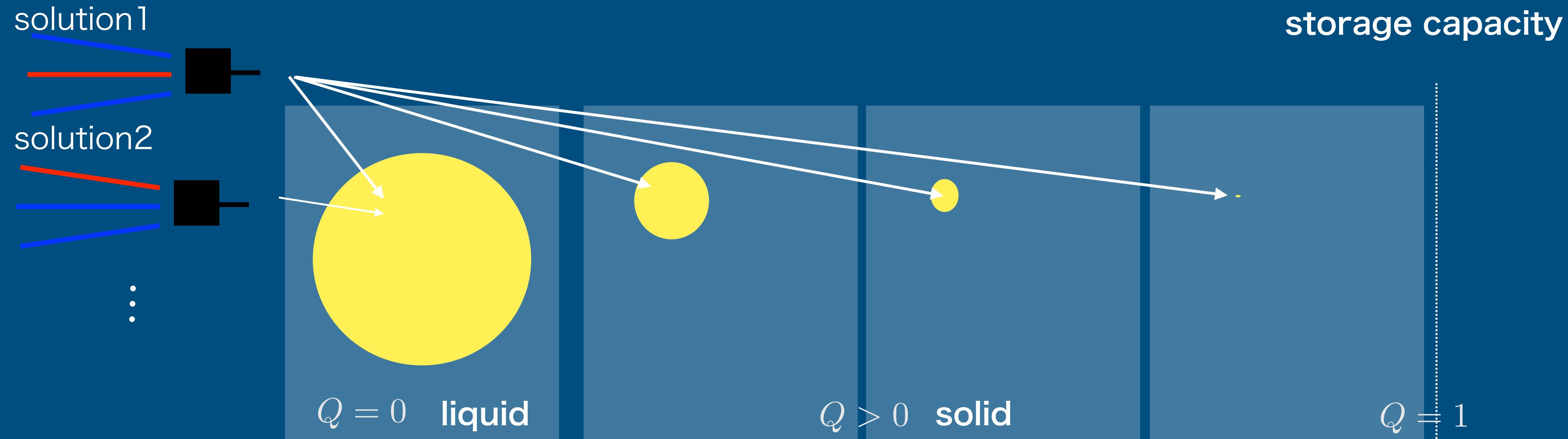
“Hardcore” constraint

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

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

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

Gardner's volume: design space of a perceptron

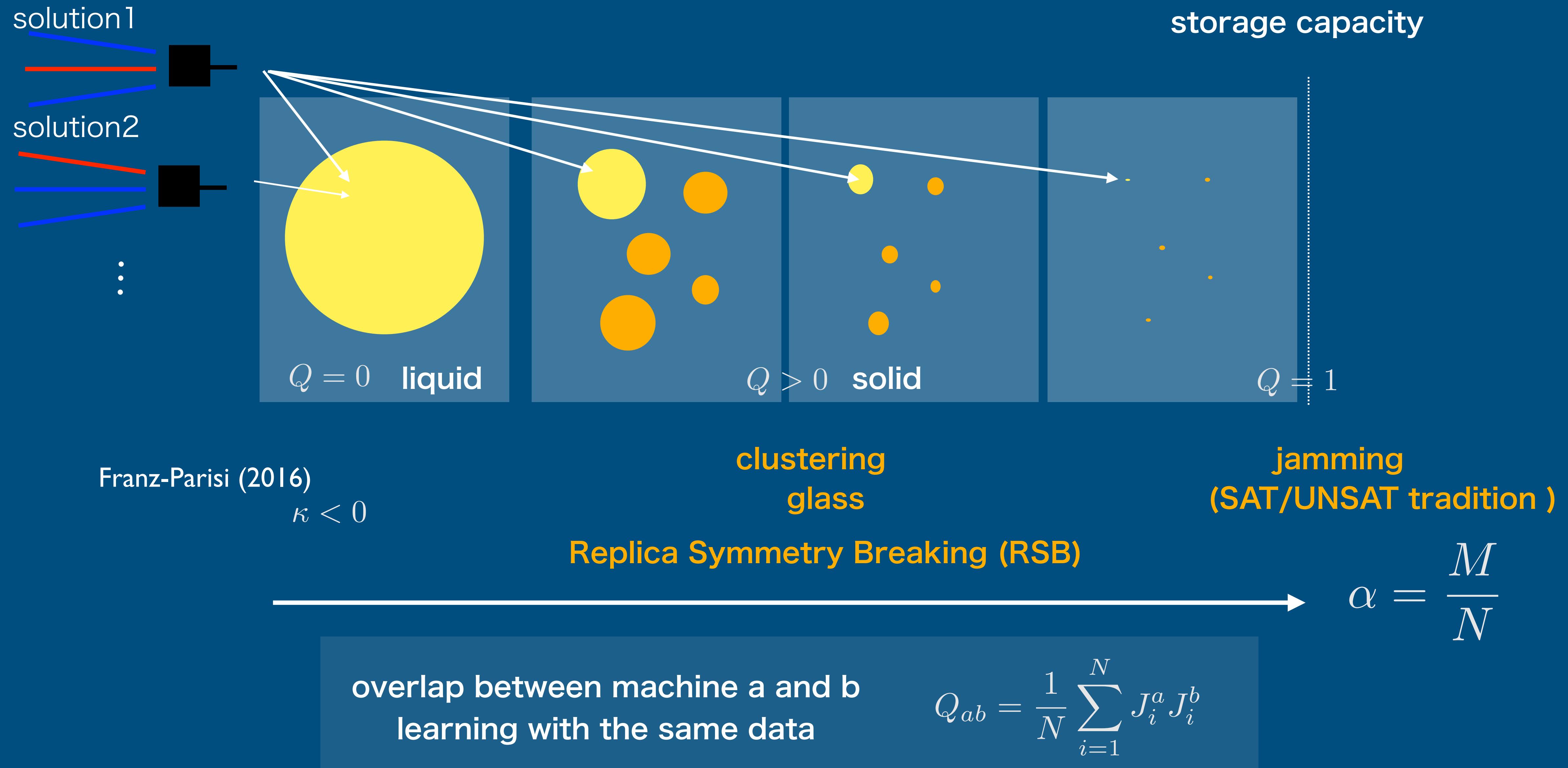


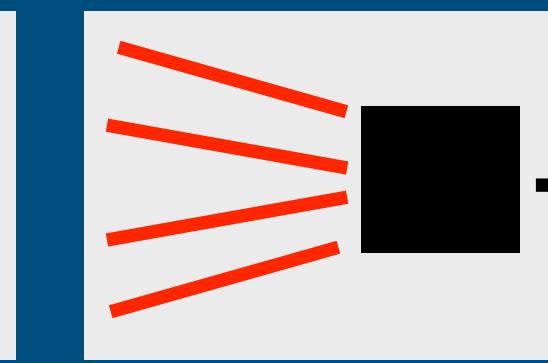
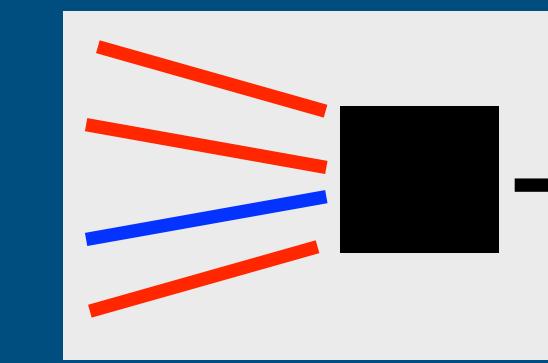
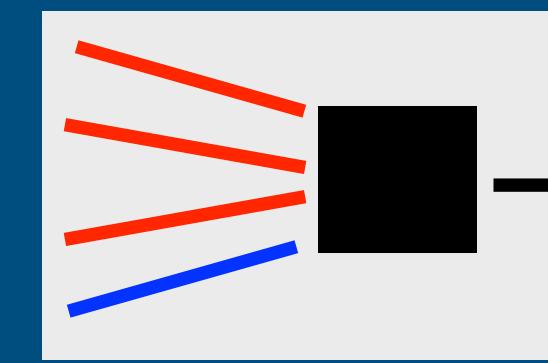
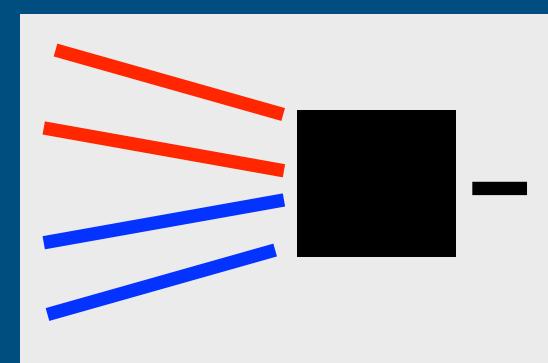
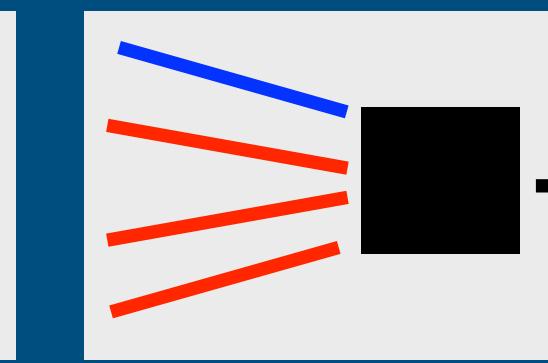
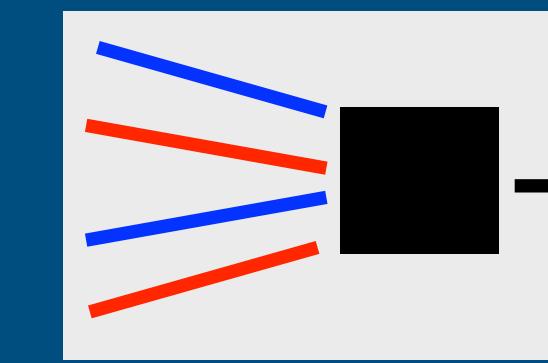
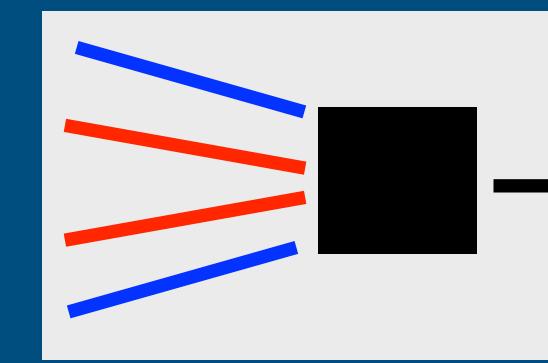
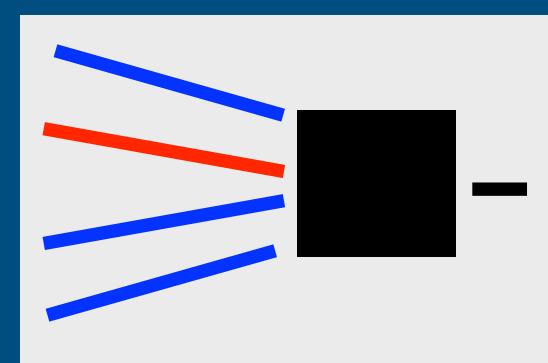
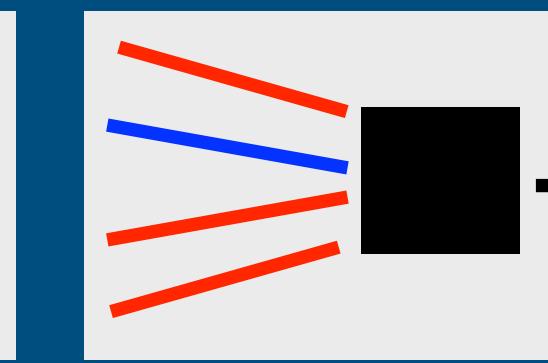
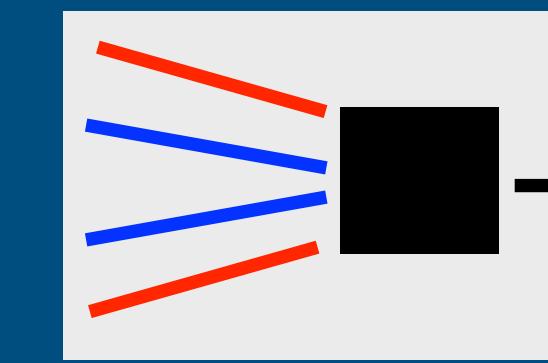
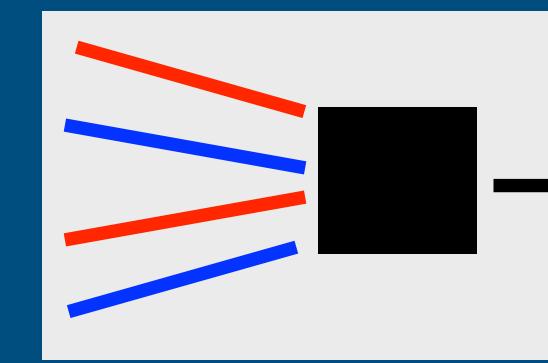
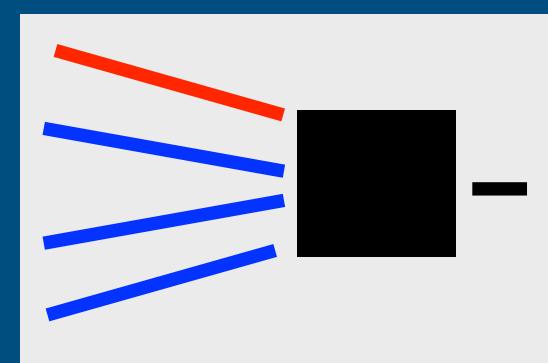
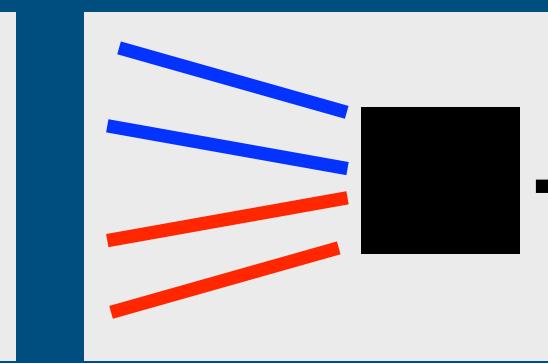
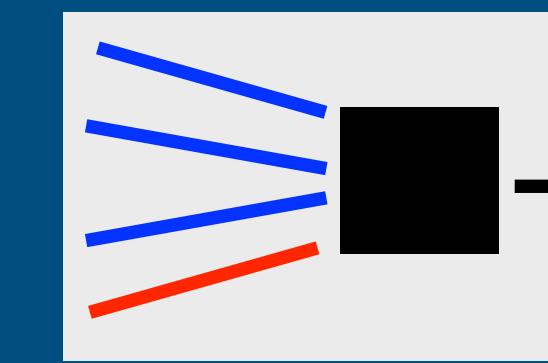
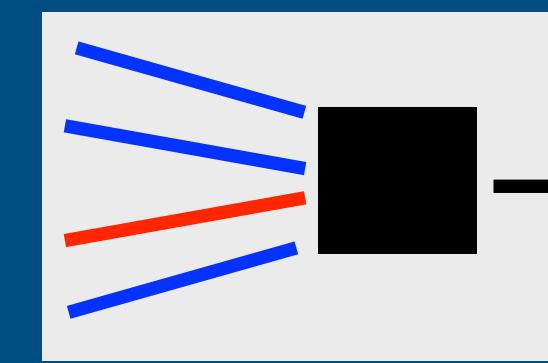
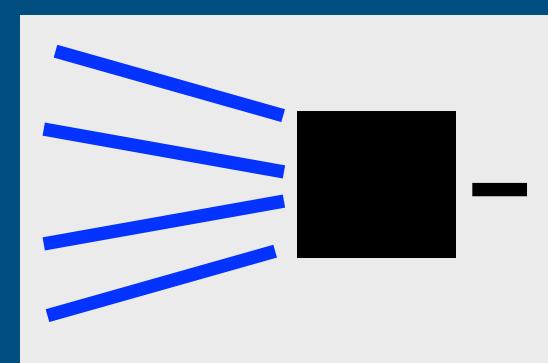
$$\alpha = \frac{M}{N}$$

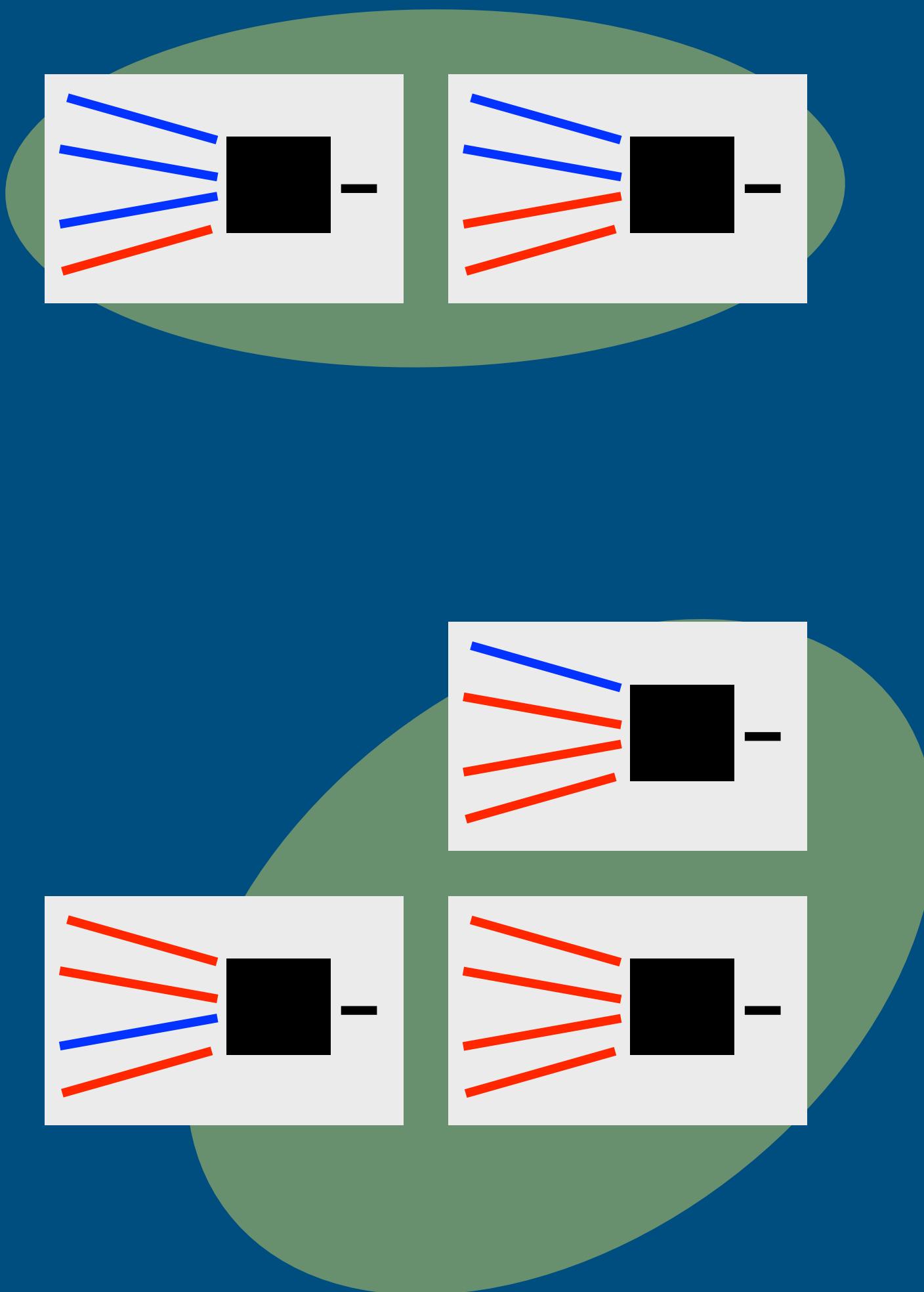
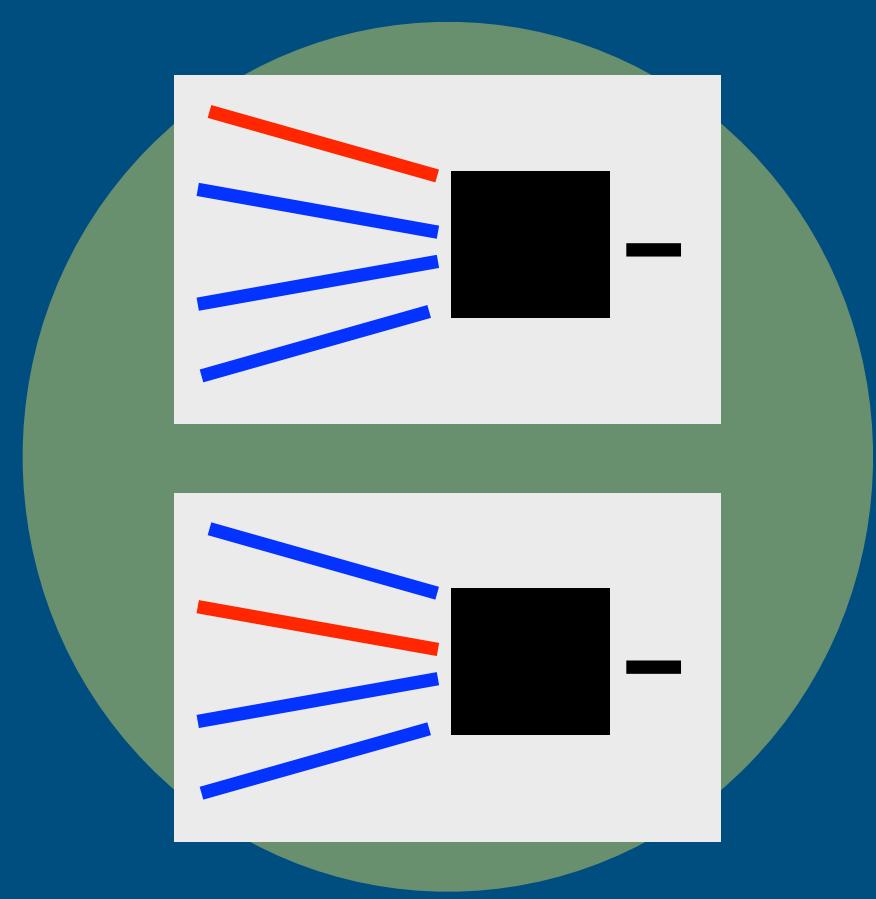
overlap between machine a and b
learning with the same data

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

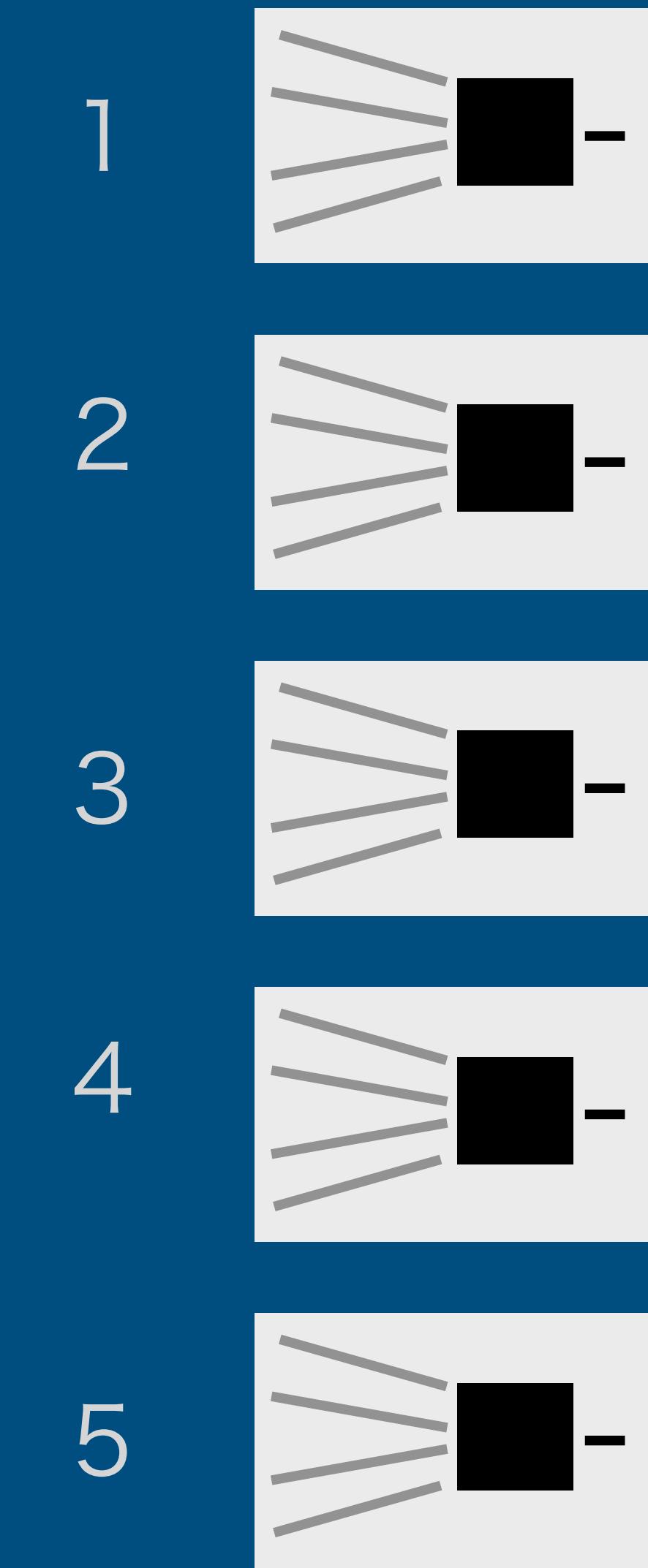
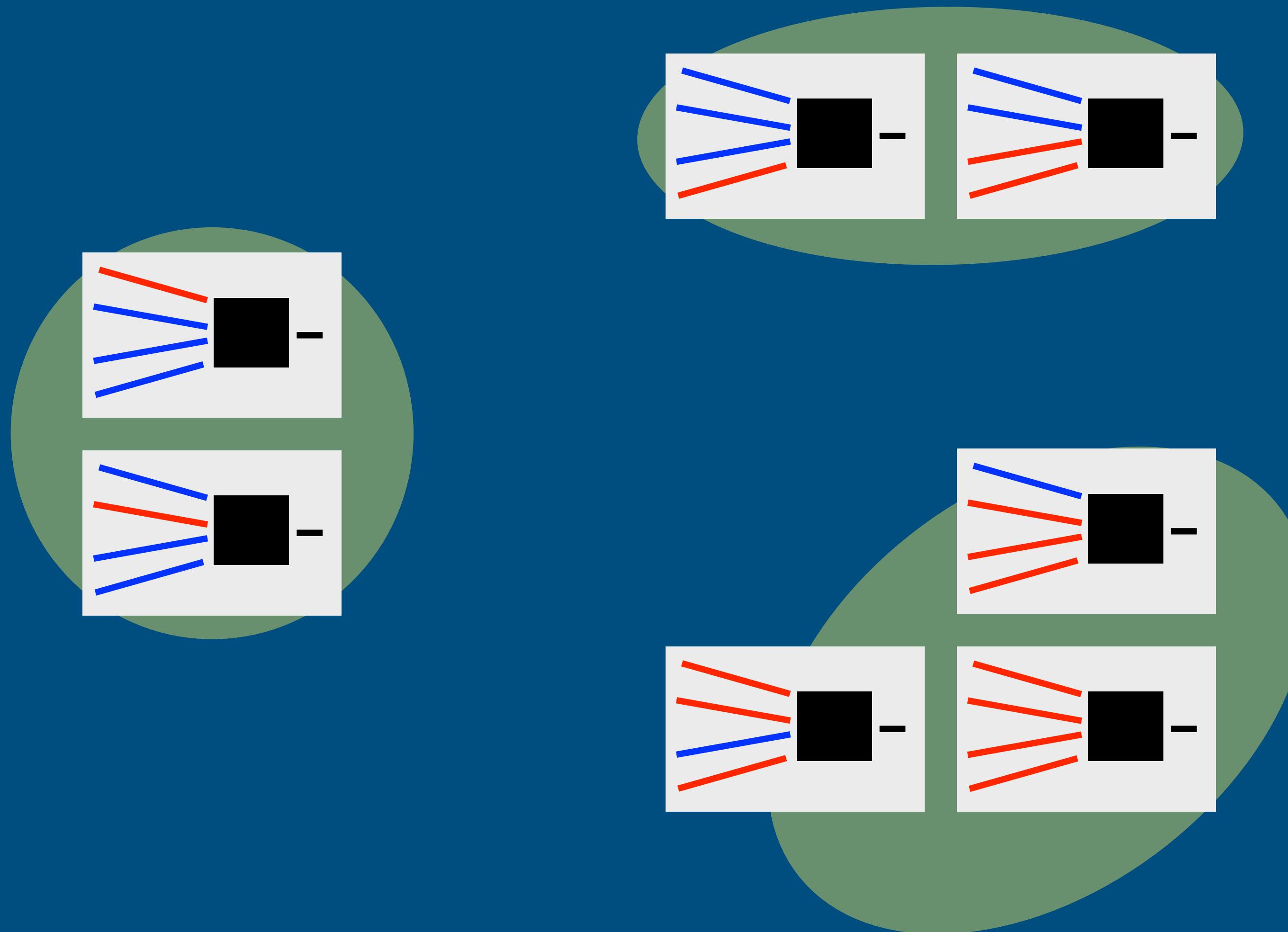
Gardner's volume: design space of a perceptron



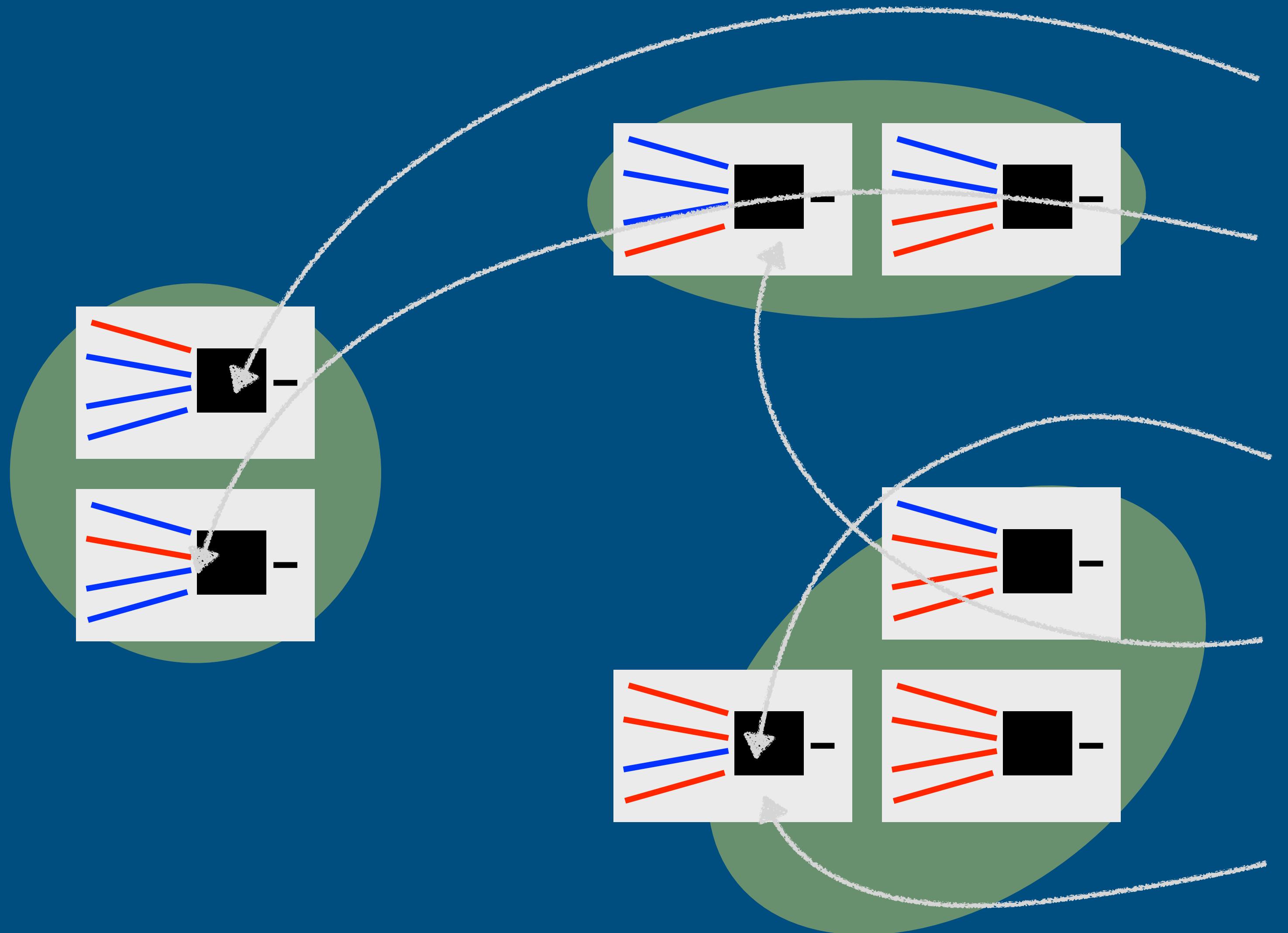




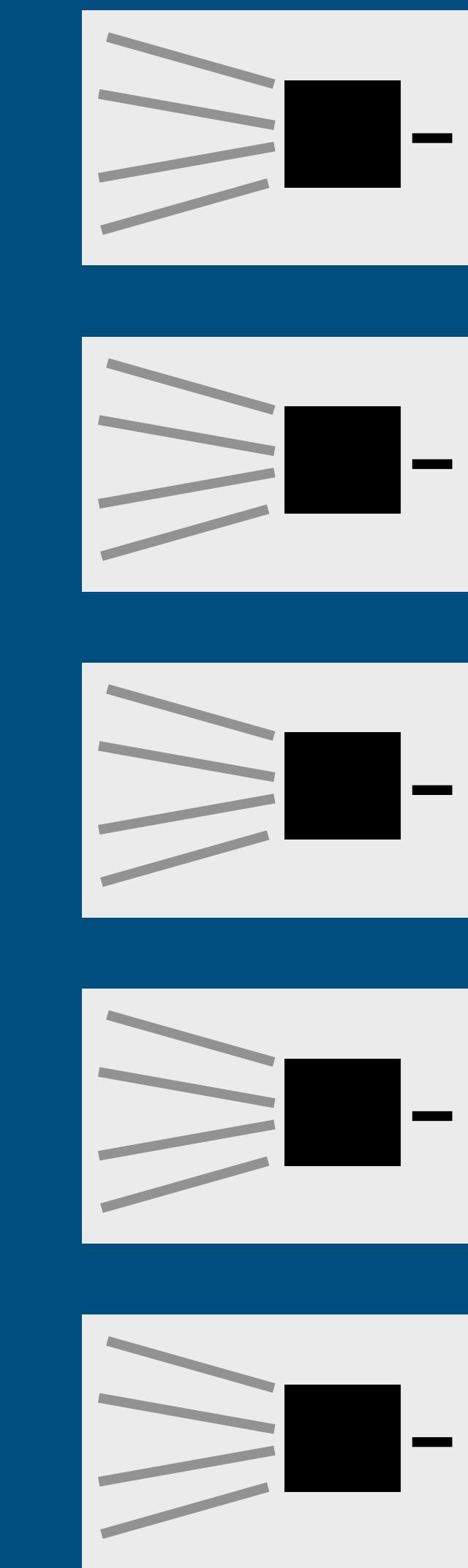
replica=machines learning in parallel
with the same training data



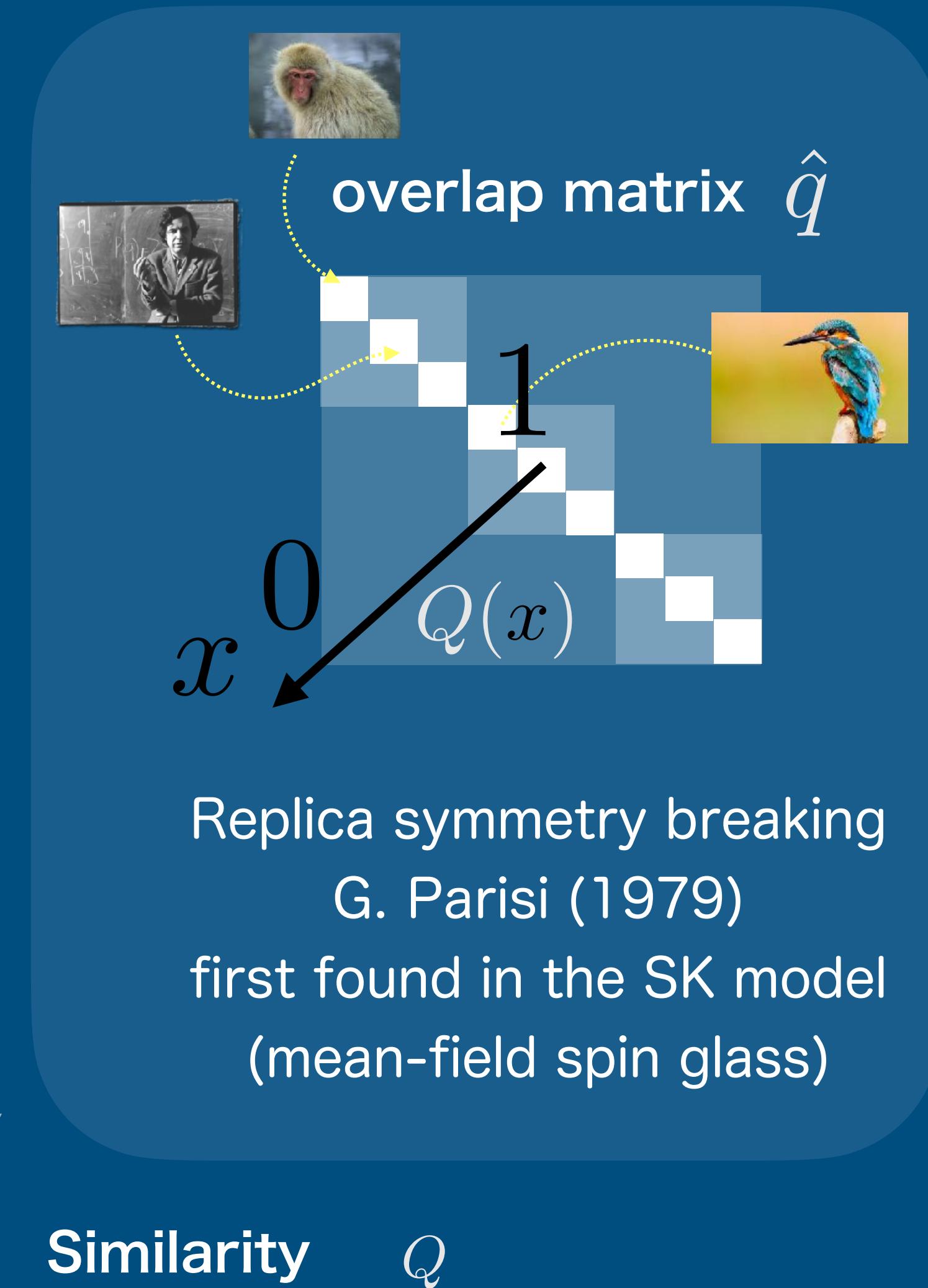
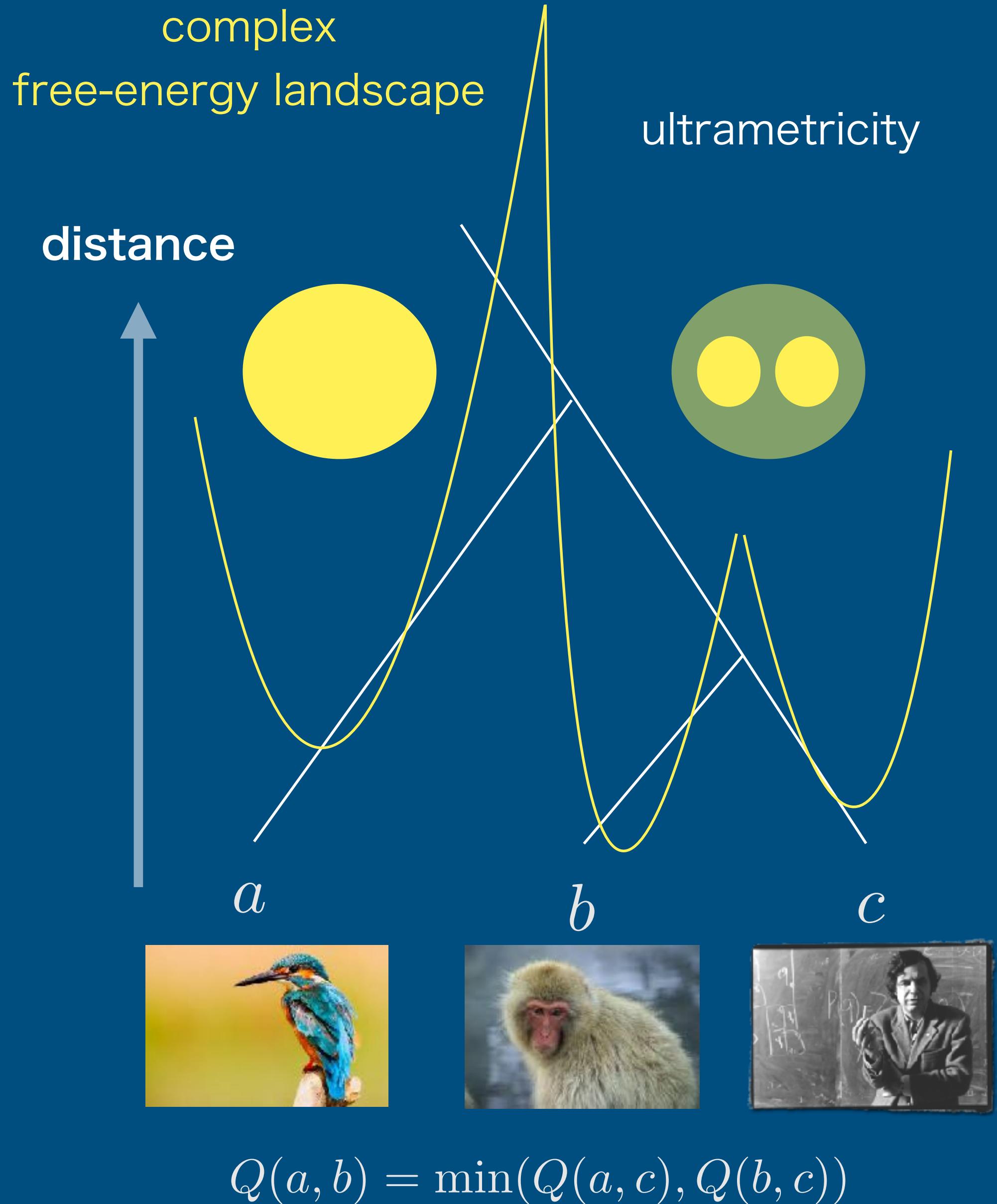
**Replica Symmetry Breaking =
Replica “Permutation” Symmetry Breaking**



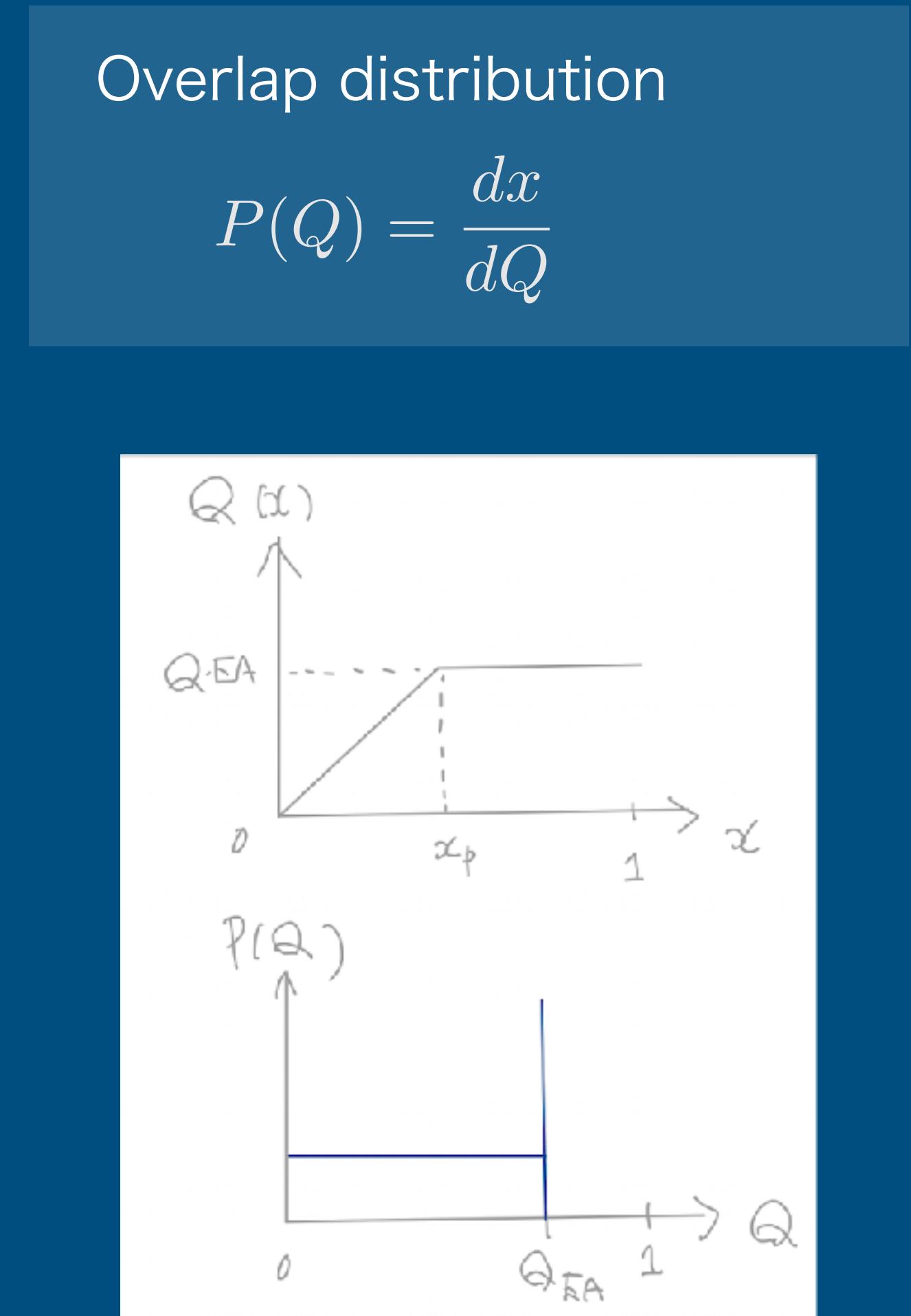
replica=machines learning in parallel
with the same training data



Hierarchical Replica (permutation) symmetry breaking and ultrametricity

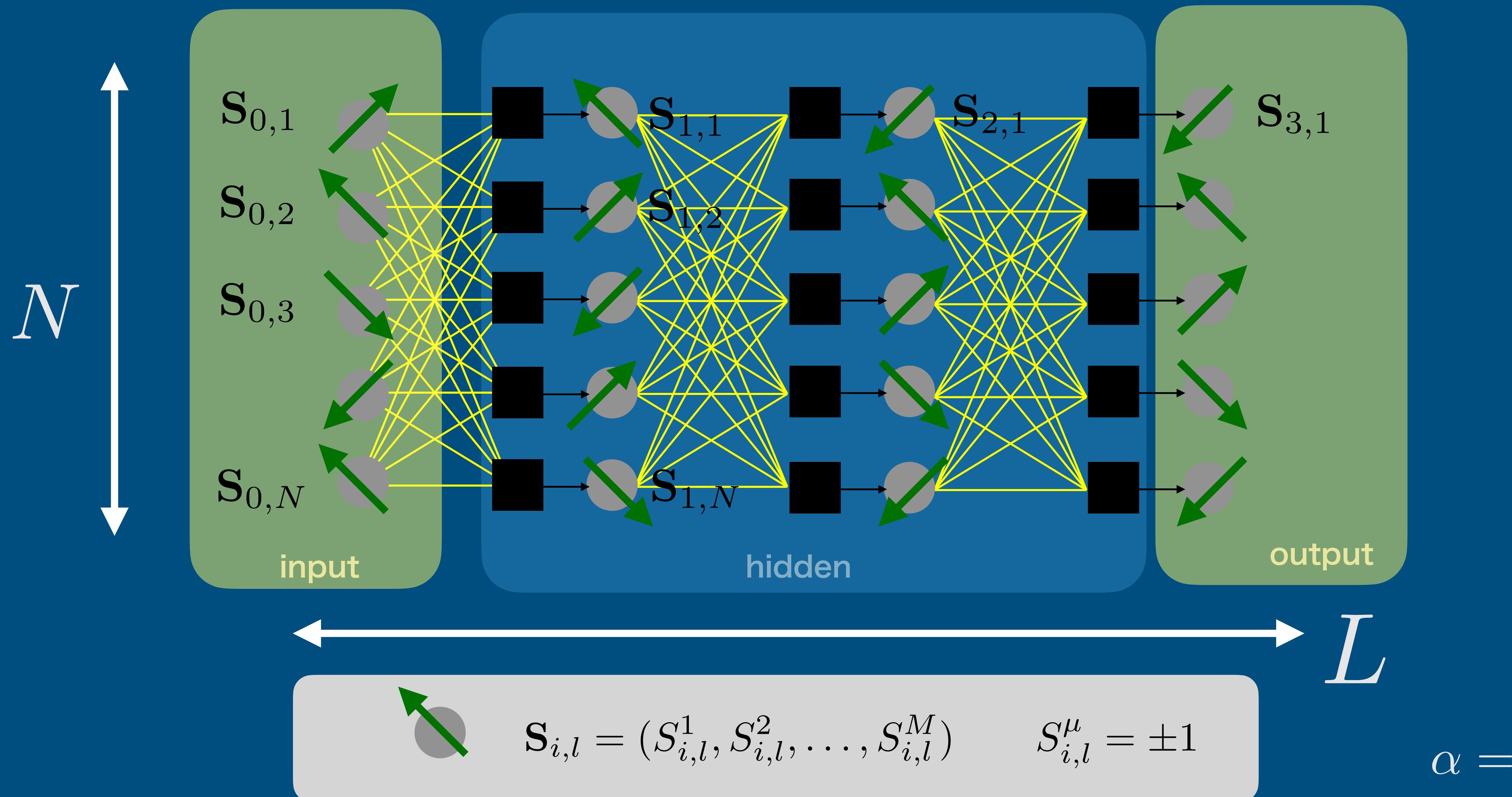


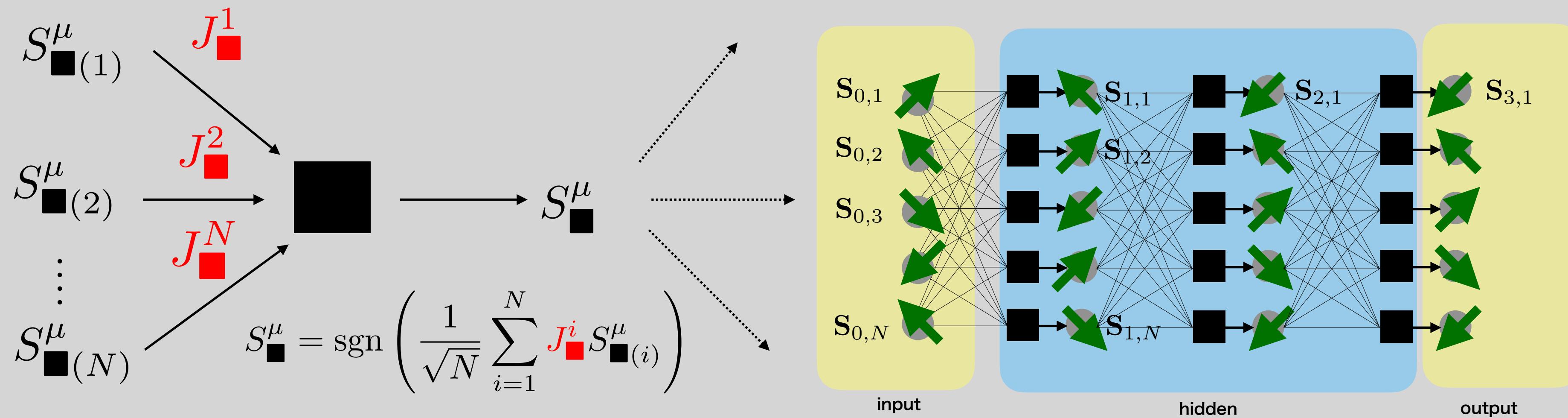
$$Q_{ab} = \frac{1}{N} \sum_{i=1}^N J_i^a J_i^b$$



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

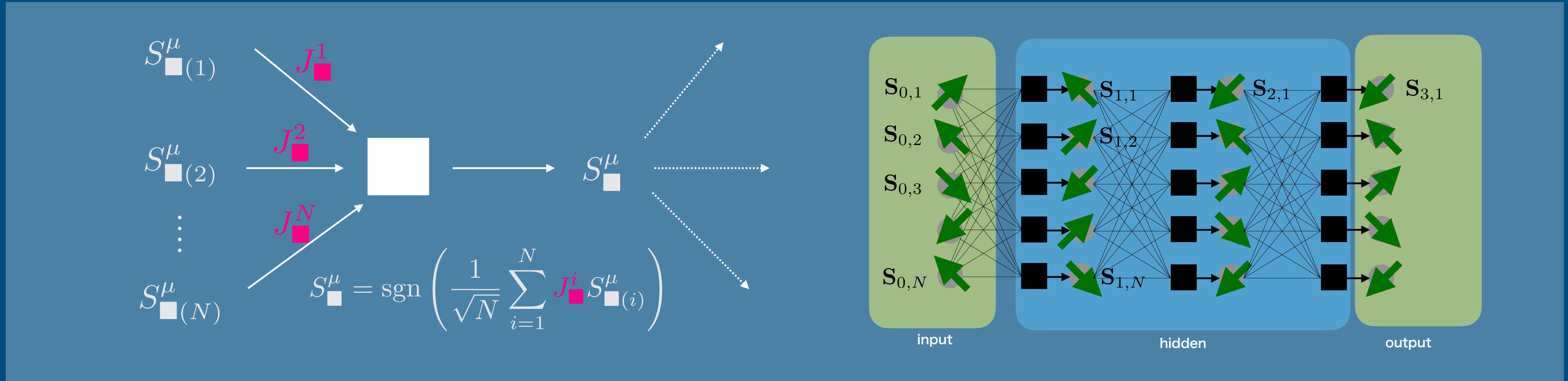
desired output ↓

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

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

Hamiltonian with
“short-ranged” interactions

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

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

“Hardcore” constraint

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

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

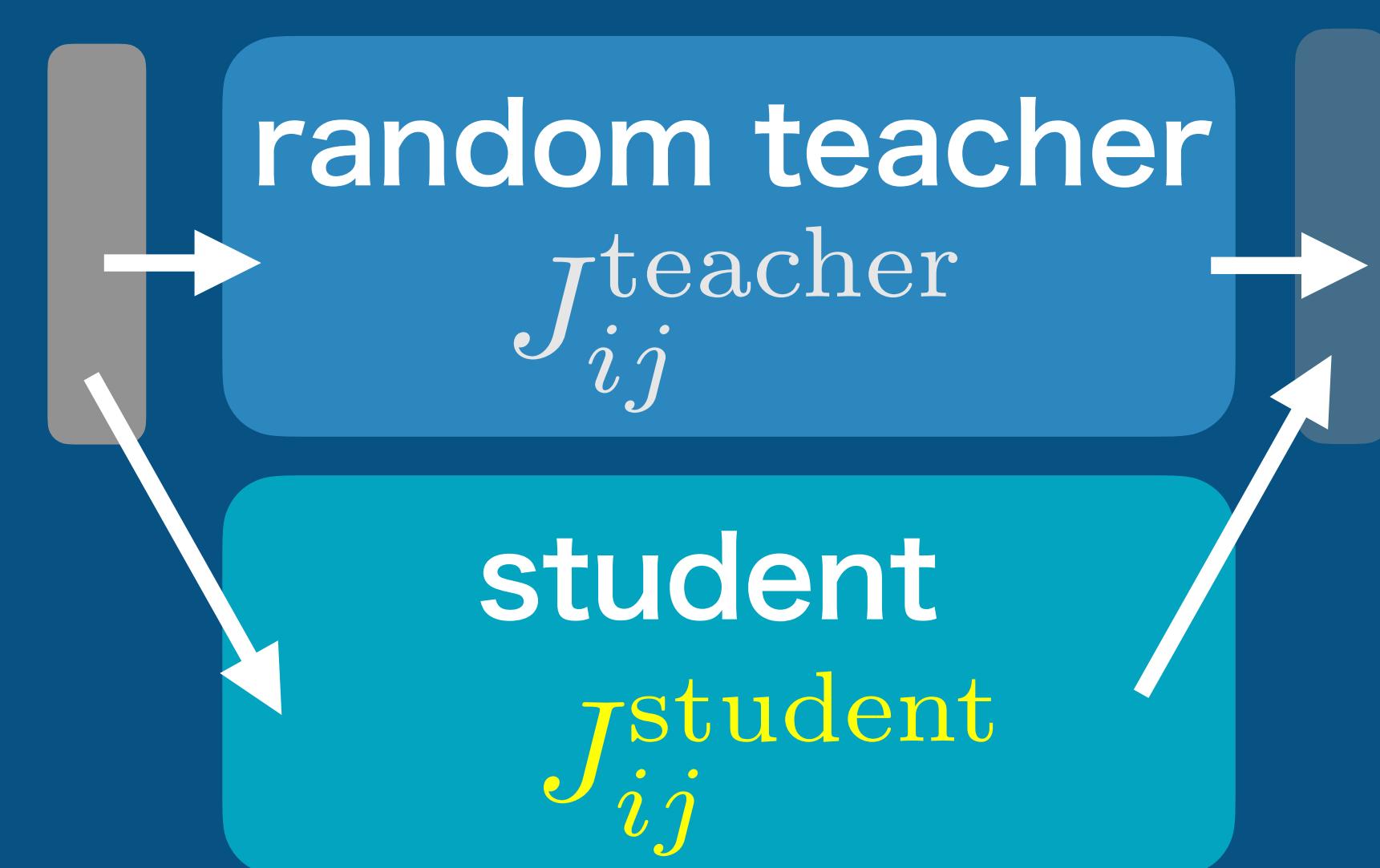
random
input



scenario (1)

random
output

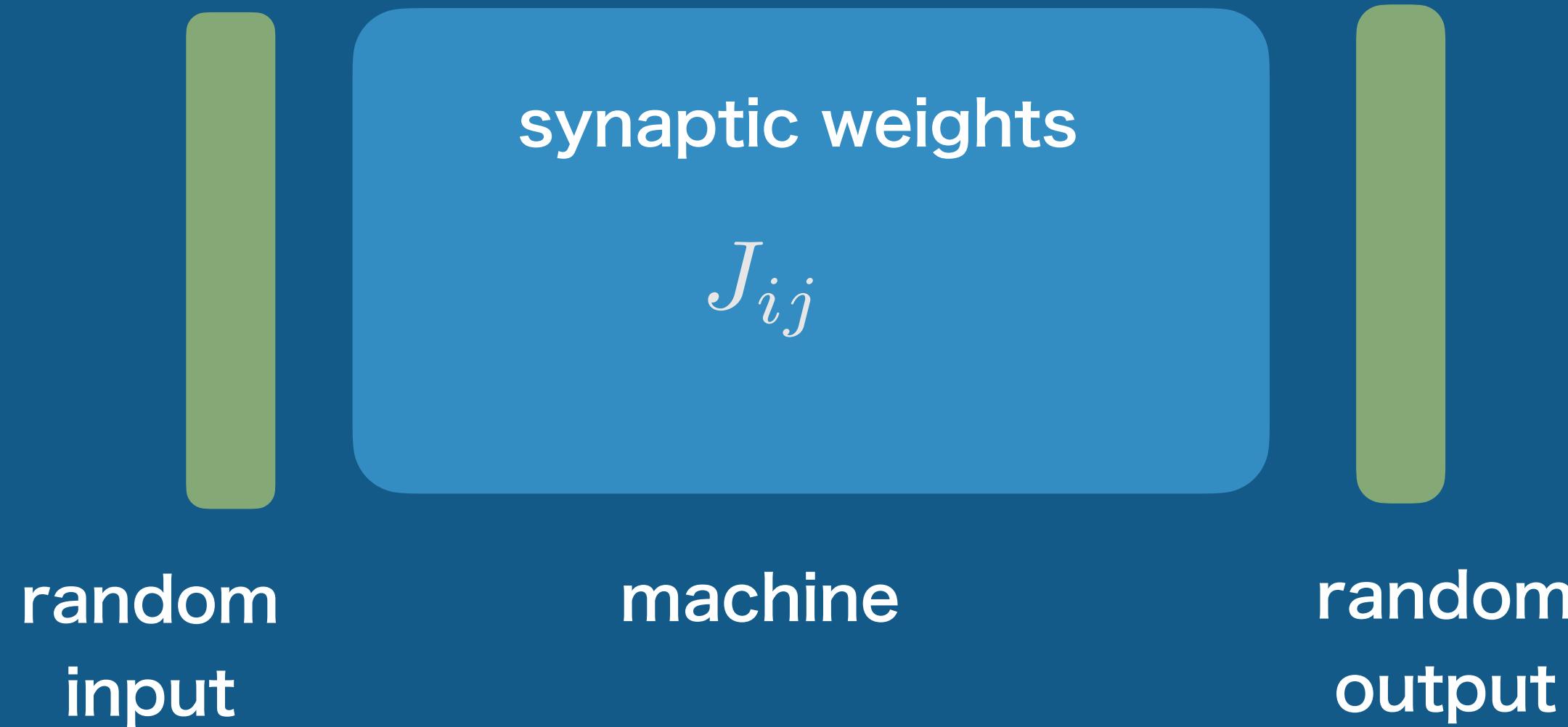
Random
input



scenario (2)

Scenario (I) Random inputs/random outputs

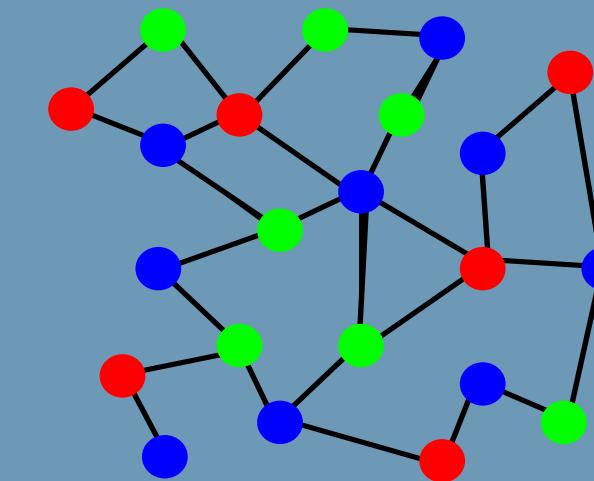
- a constraint satisfaction problem
- glass/jamming physics



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

“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

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

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

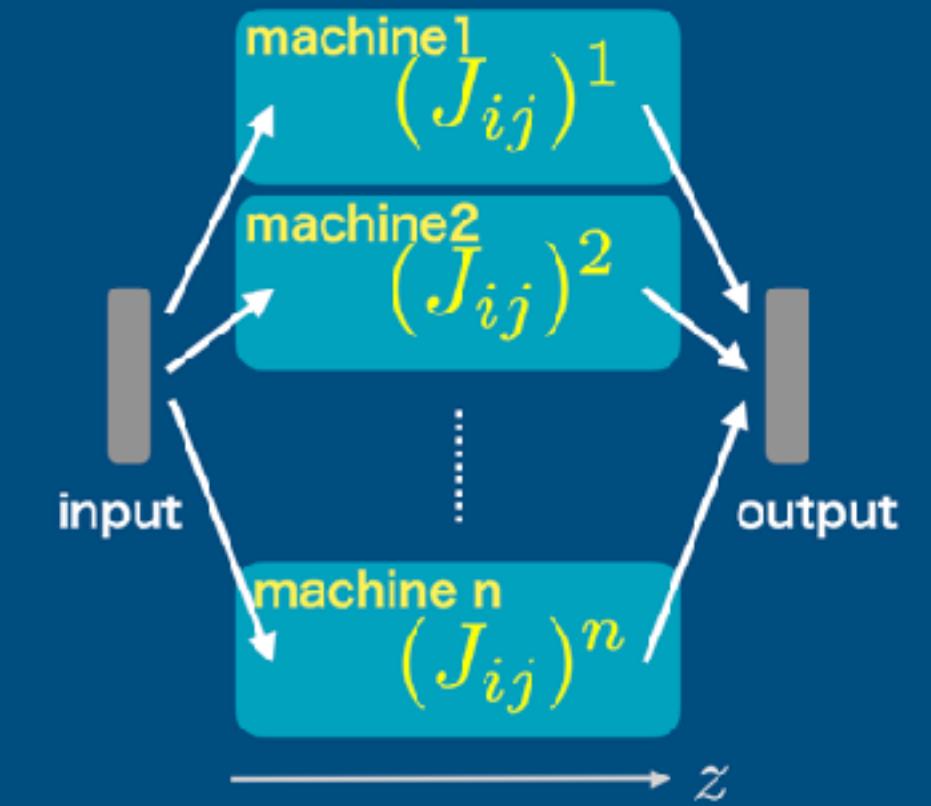
$$\frac{-\beta \overline{F(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} = \frac{\partial_n \overline{V^n(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} \Big|_{n=0} = \partial_n S_n[\{\hat{Q}(l), \hat{q}(l)\}] \Big|_{n=0}$$

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

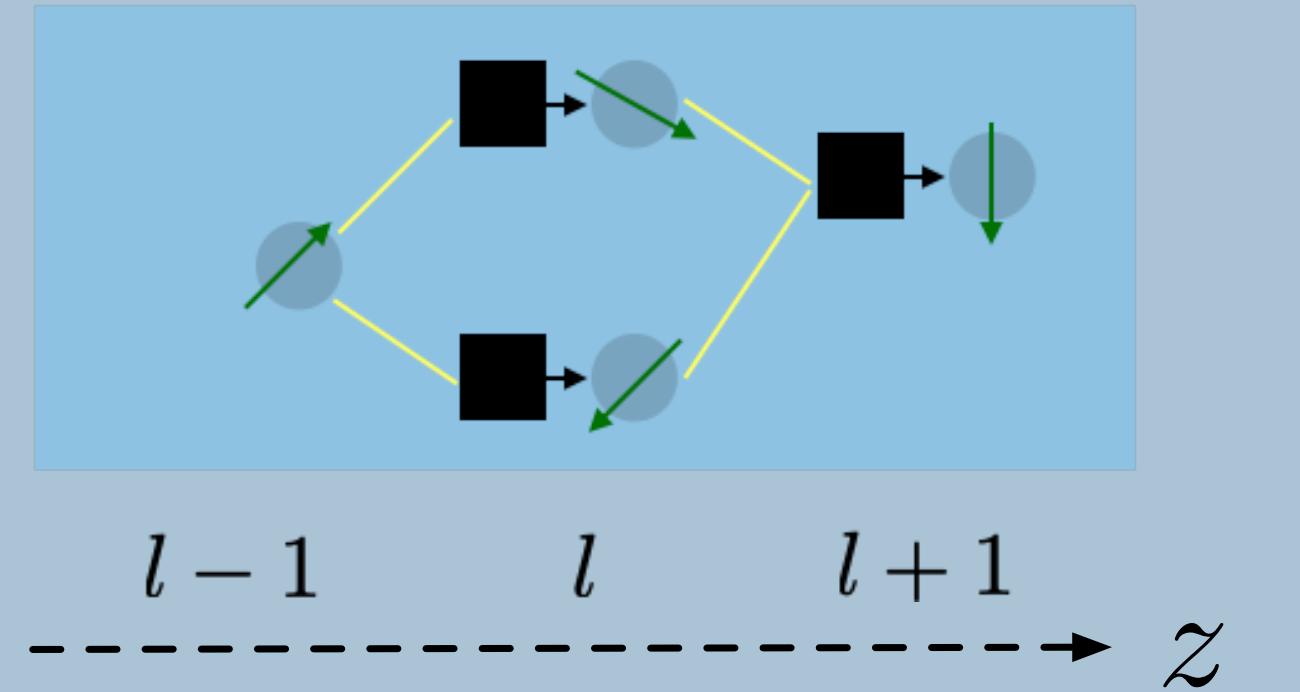
$$\alpha = \frac{M}{N}$$

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

replicas: machines learning in parallel with the same data



Loop correction

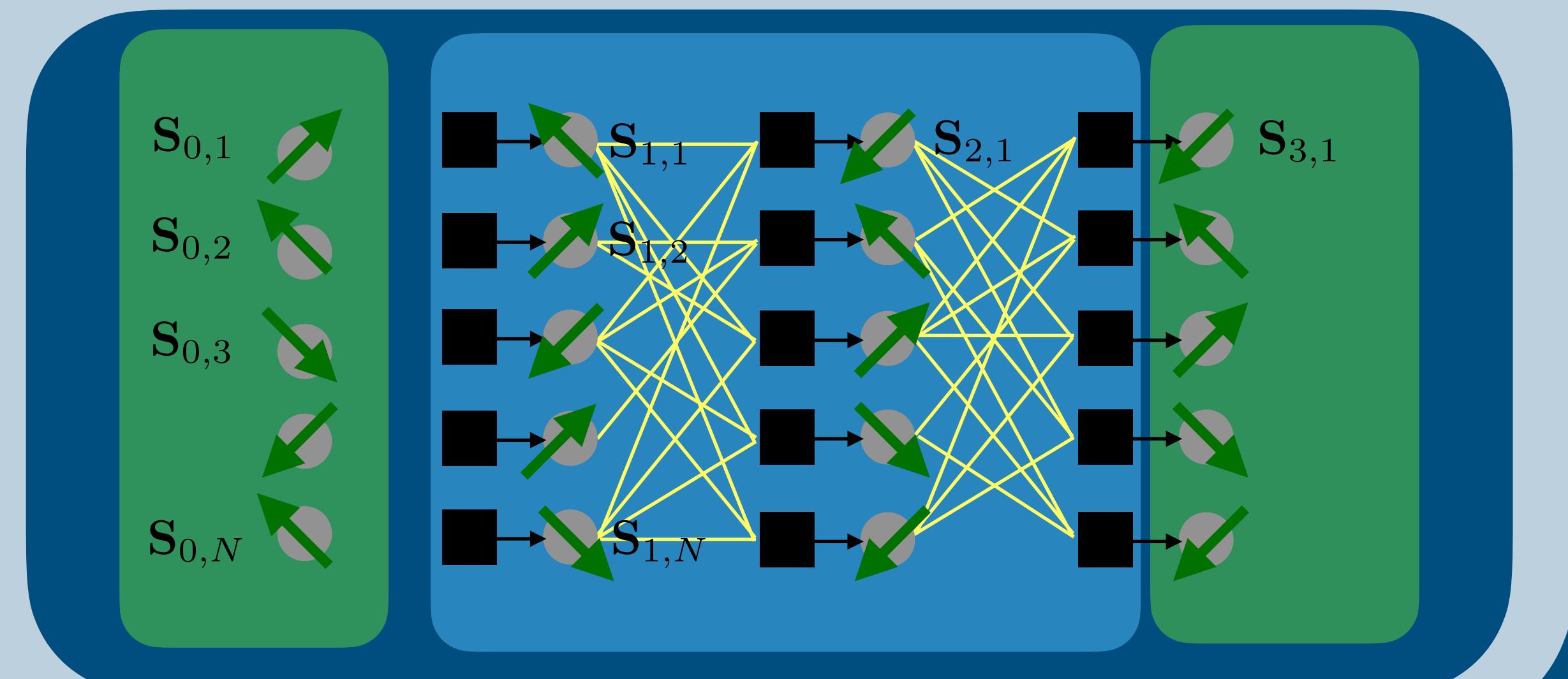


Not global but **dense** coupling - loop corrections become negligible

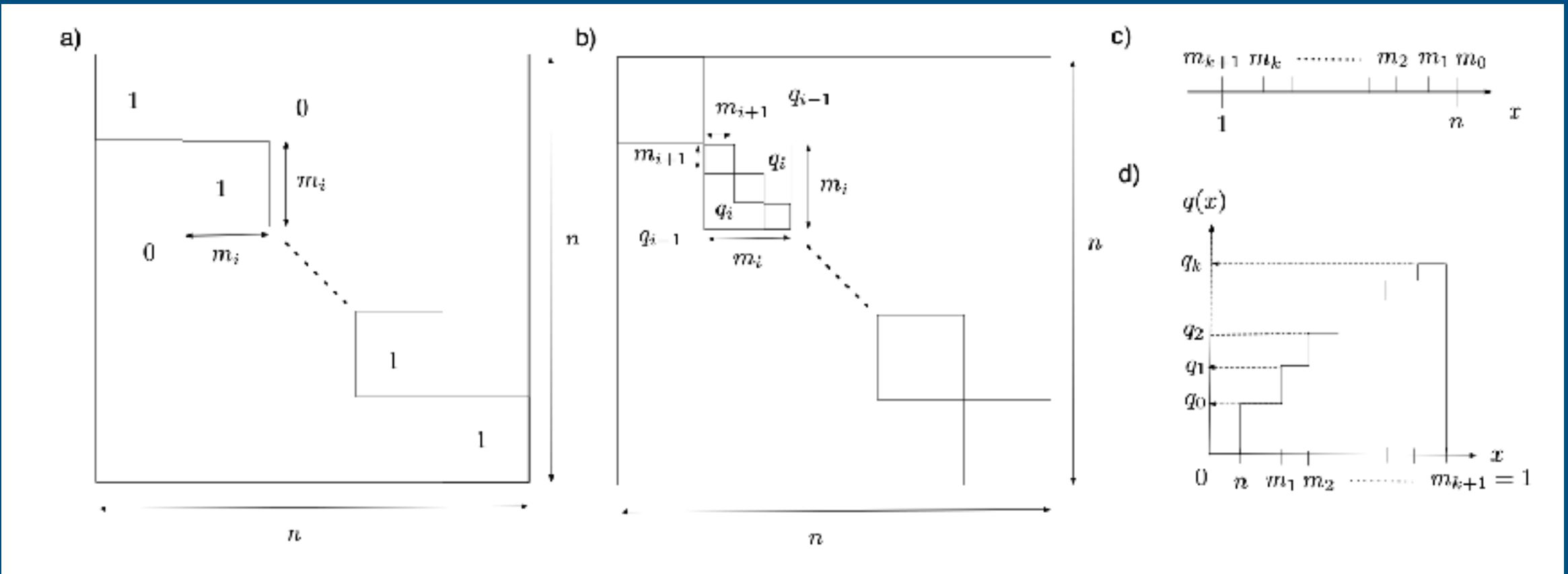
“connectivity” C

$$N \gg c \gg 1$$

$$\alpha = M/c$$



■ Parisi's RSB ansatz



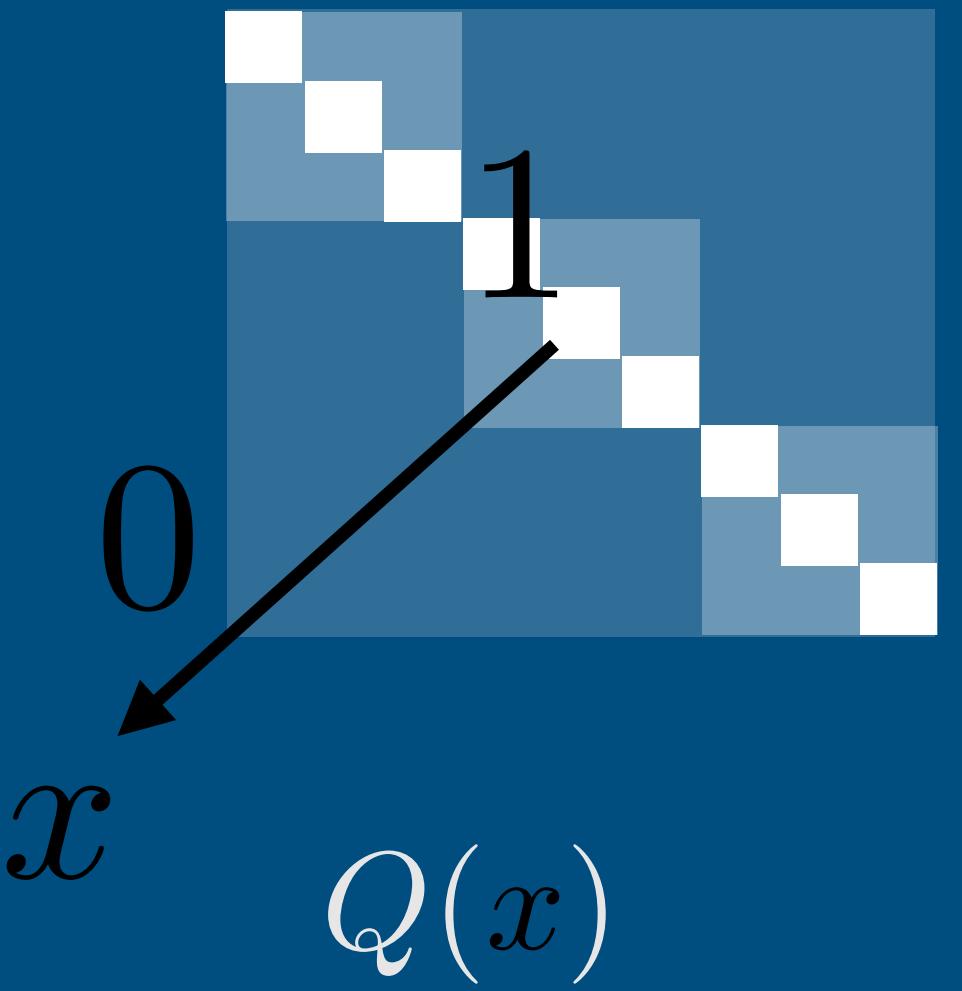
$$Q_{ab}(l) = \sum_{i=0}^{k+1} Q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L$$

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

■ Input/output boundaries

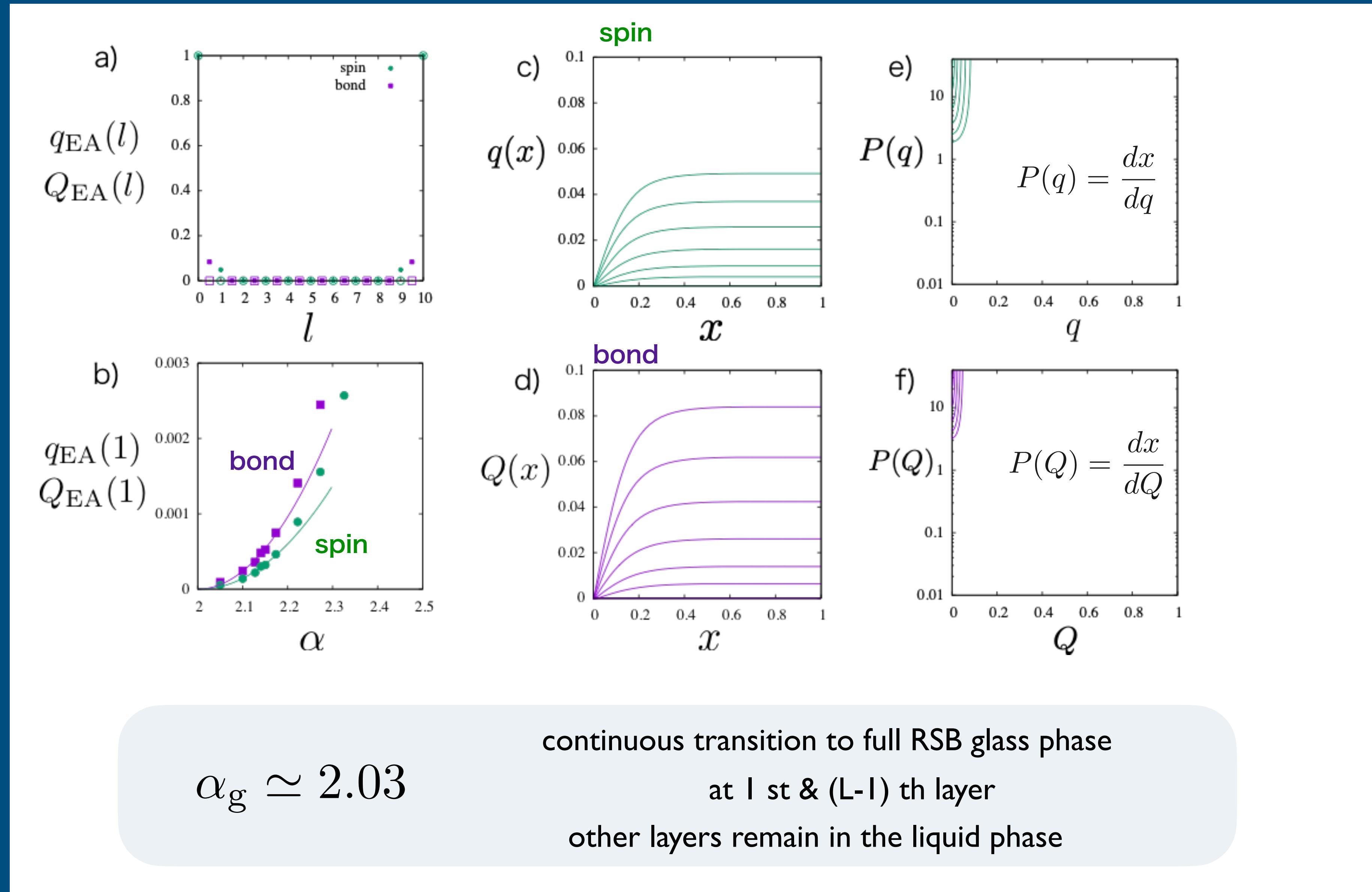
replicated machines are subjected to
the same training data

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

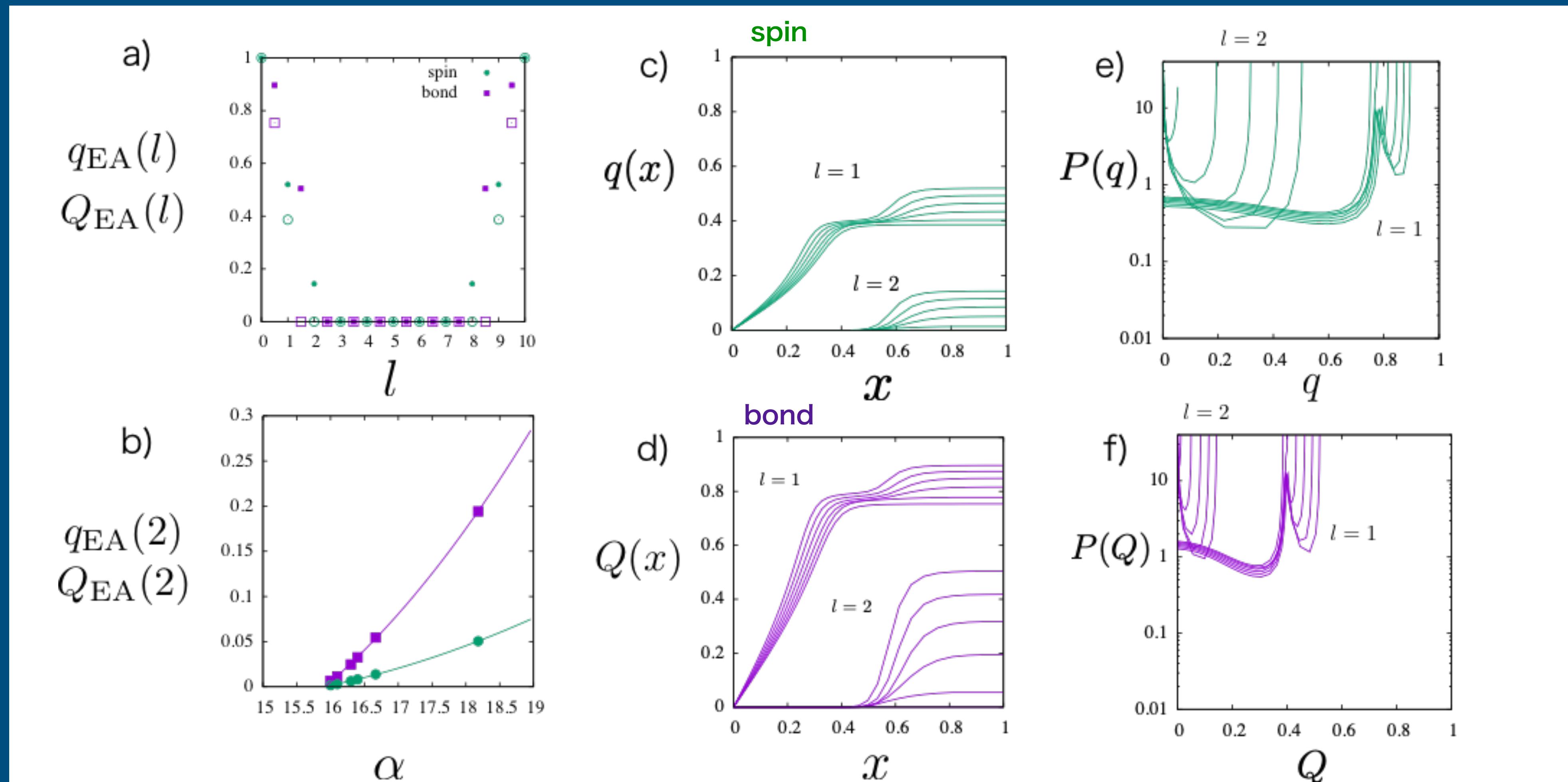


$$q(x)$$

■ 1st Glass transition



■ 2nd Glass transition

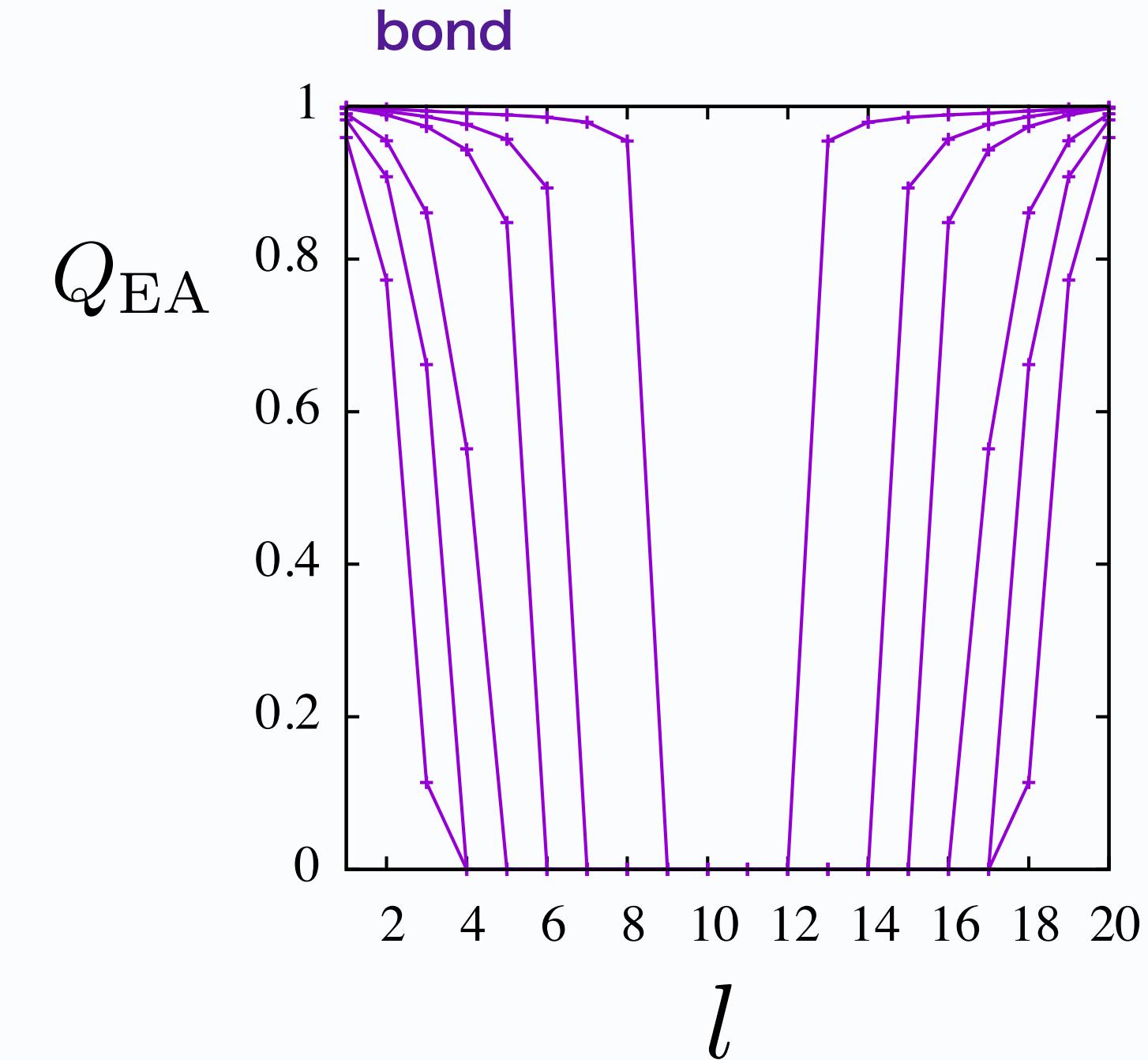
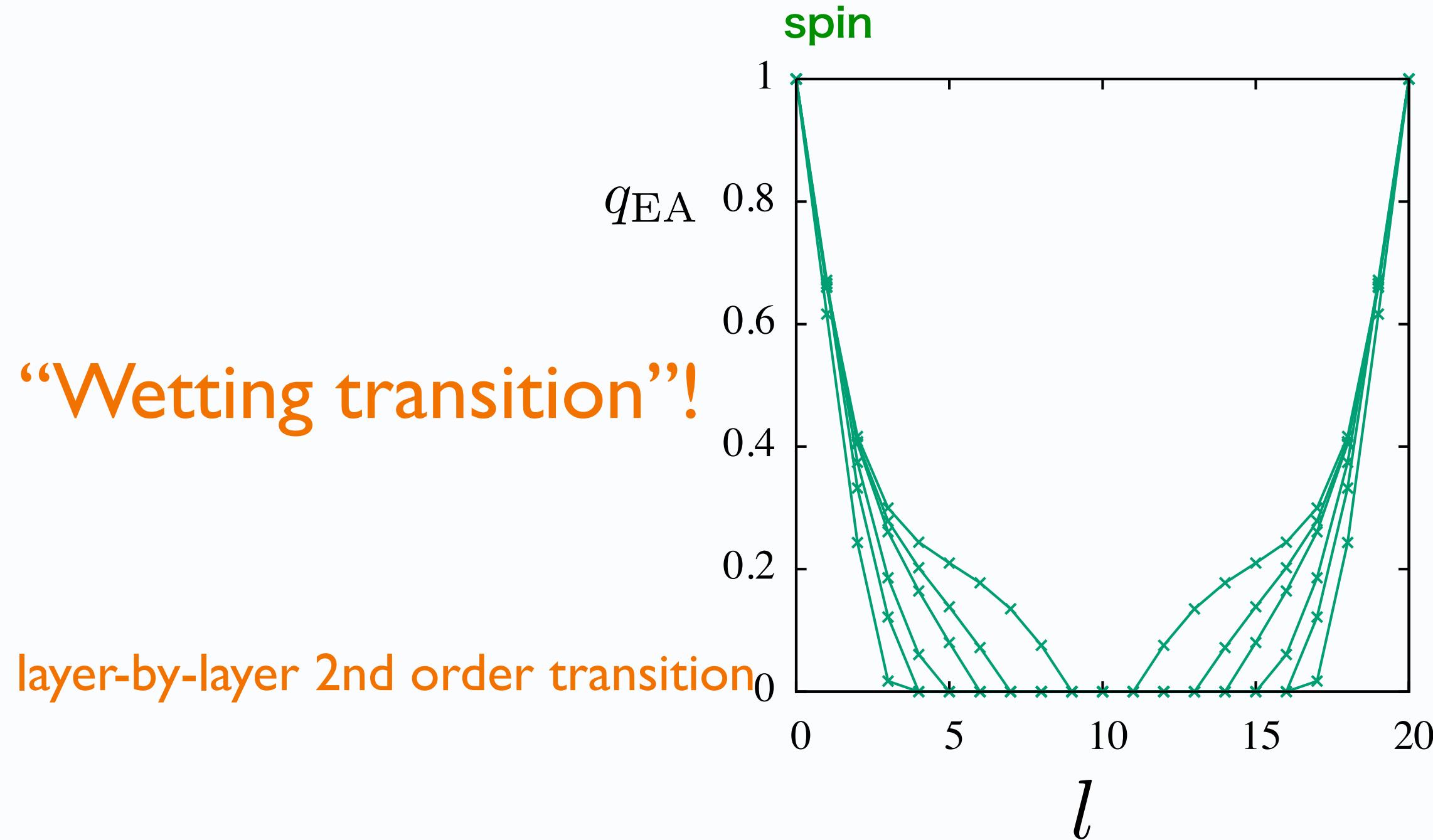


$$\alpha_g(2) \simeq 15.38$$

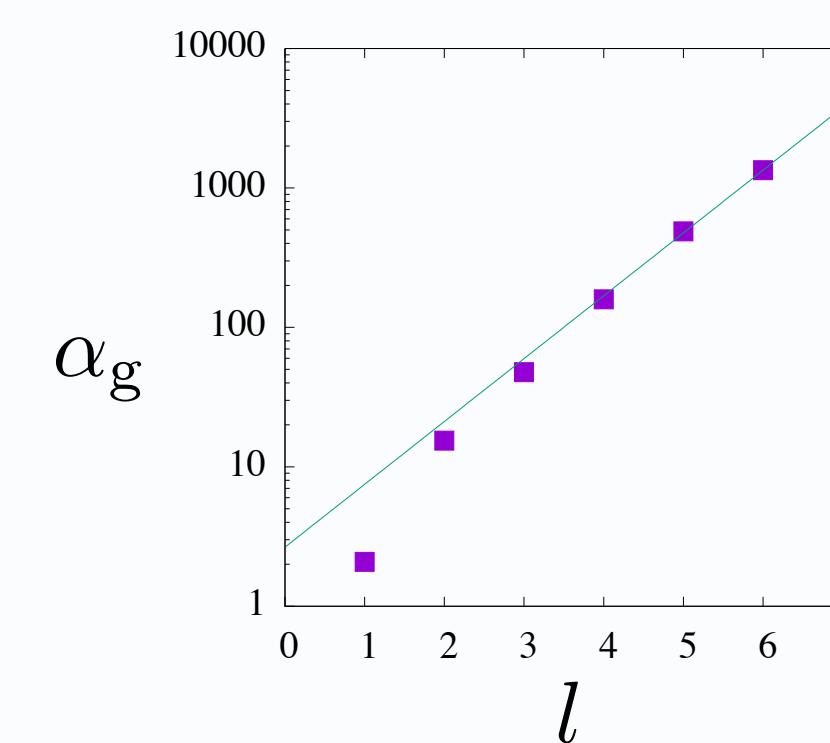
continuous transition to full RSB glass phase
at 2nd & (L-2) th layer

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

Growth of glass phase with increasing training data



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



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

this suggests storage capacity also grows fast with the depth

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

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

“penetration depth”

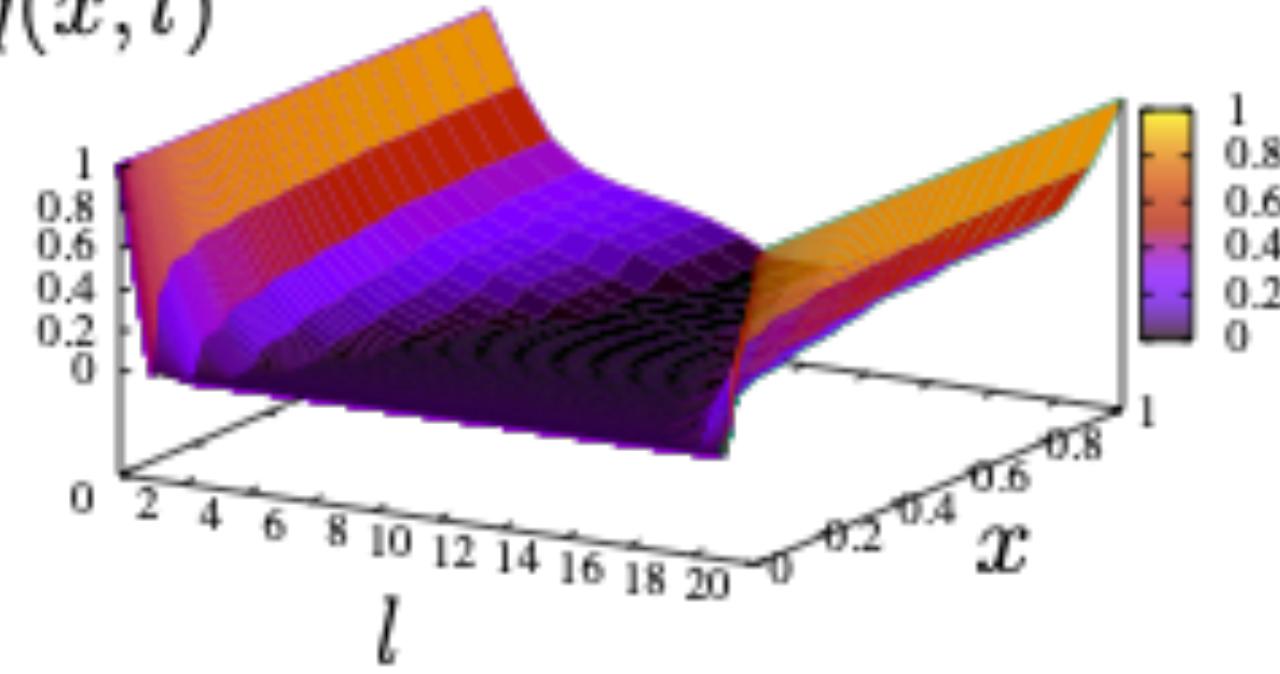
■ Space-dependent replica-symmetry breaking

$\alpha = 4000$

a)

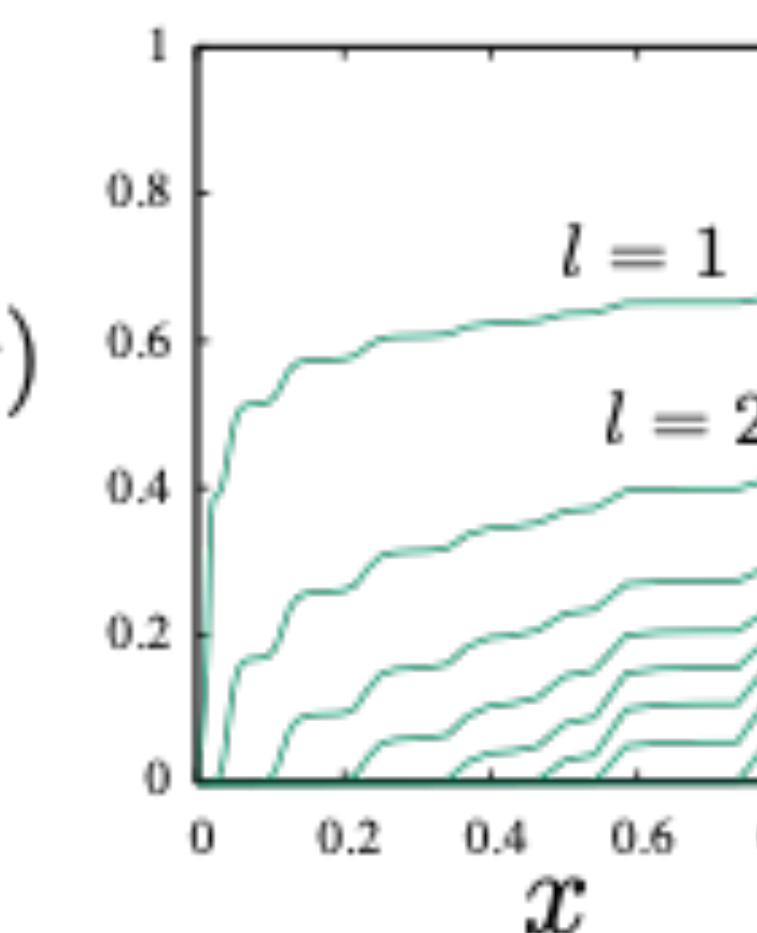
spin

$q(x, l)$

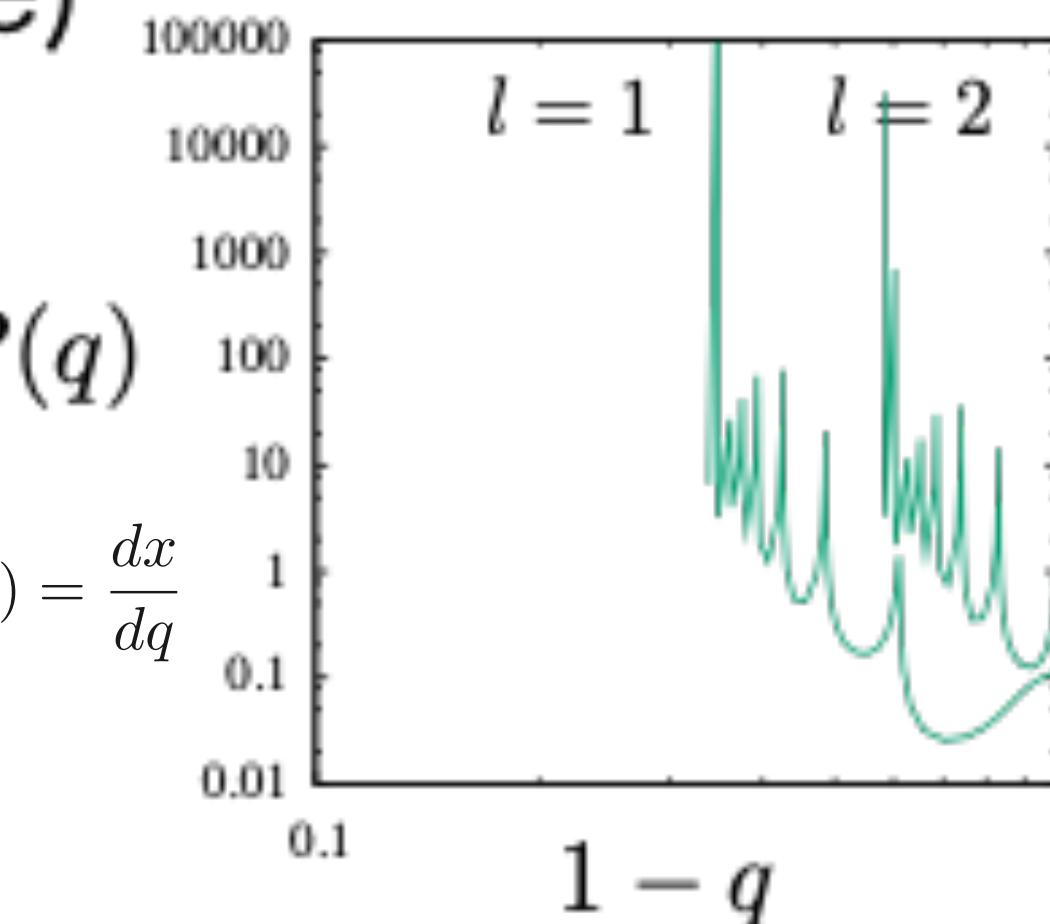


c)

$q(x)$



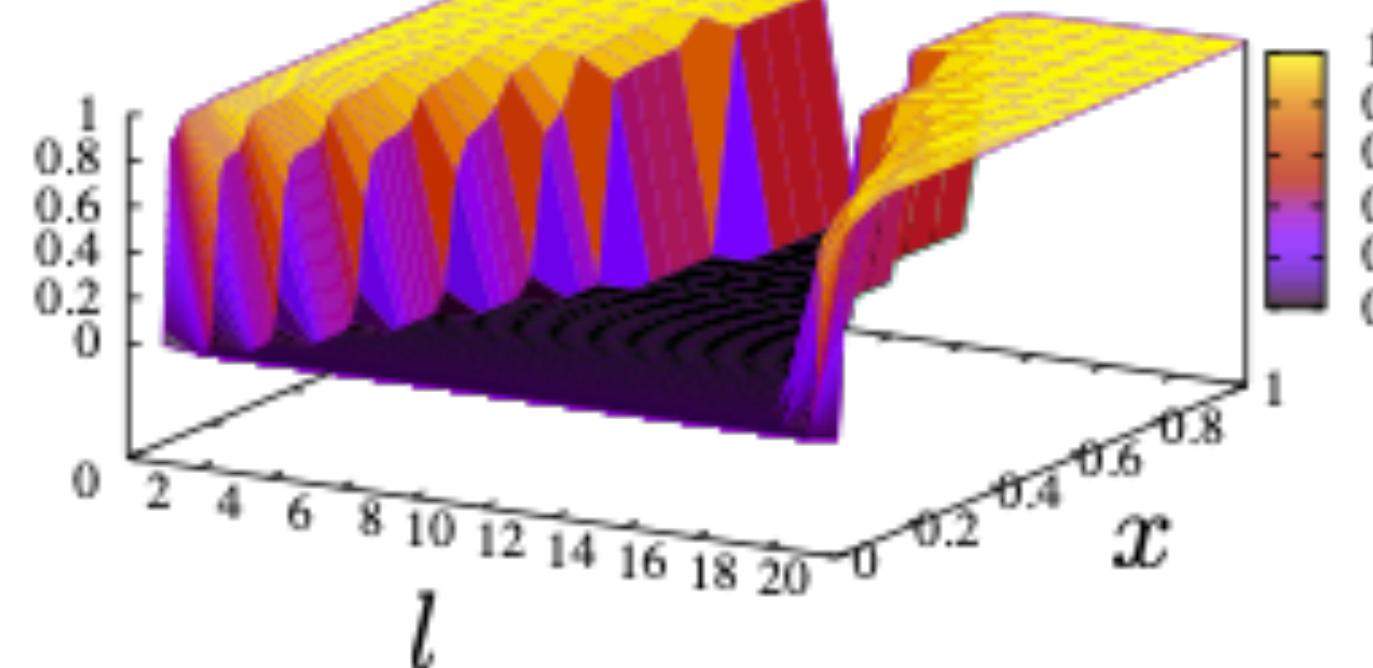
e)



b)

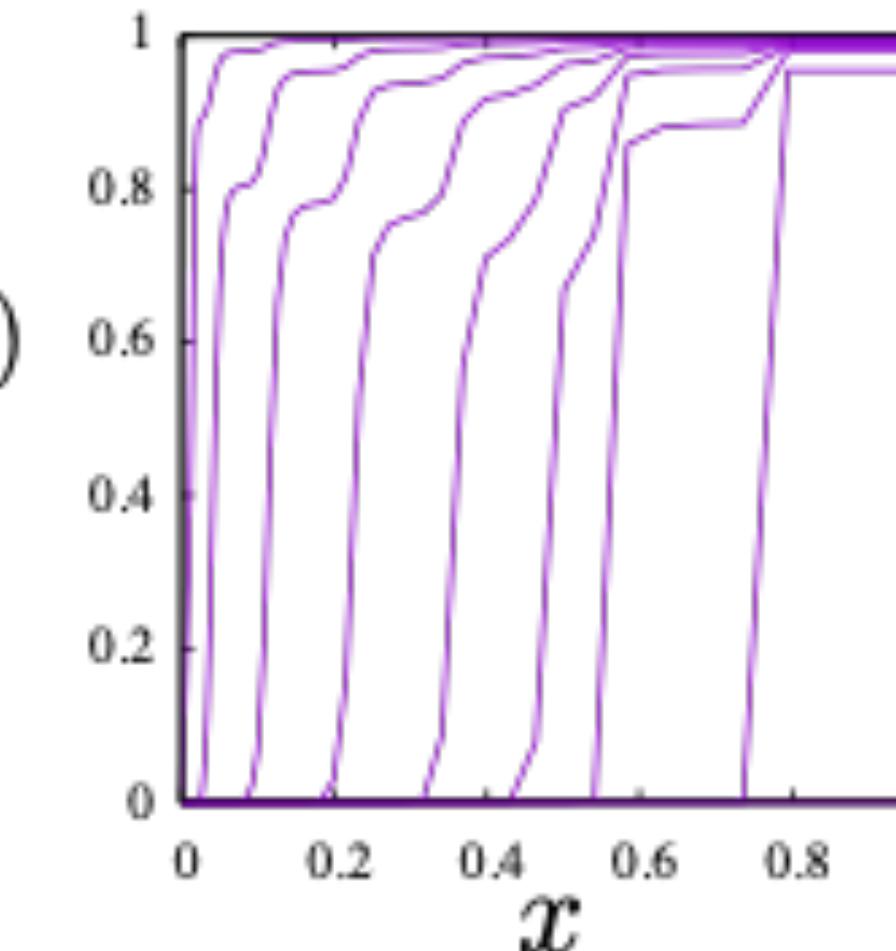
bond

$Q(x, l)$

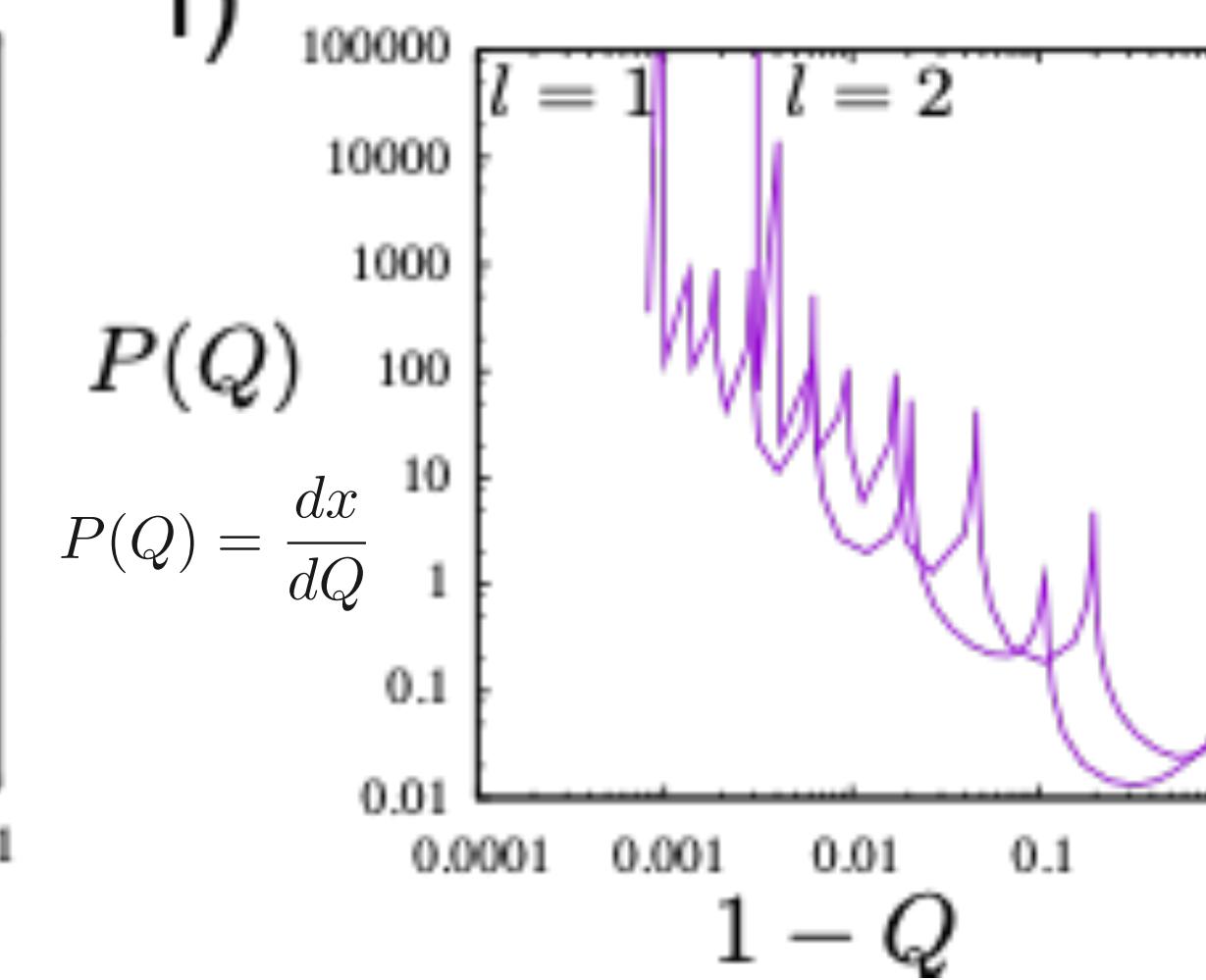


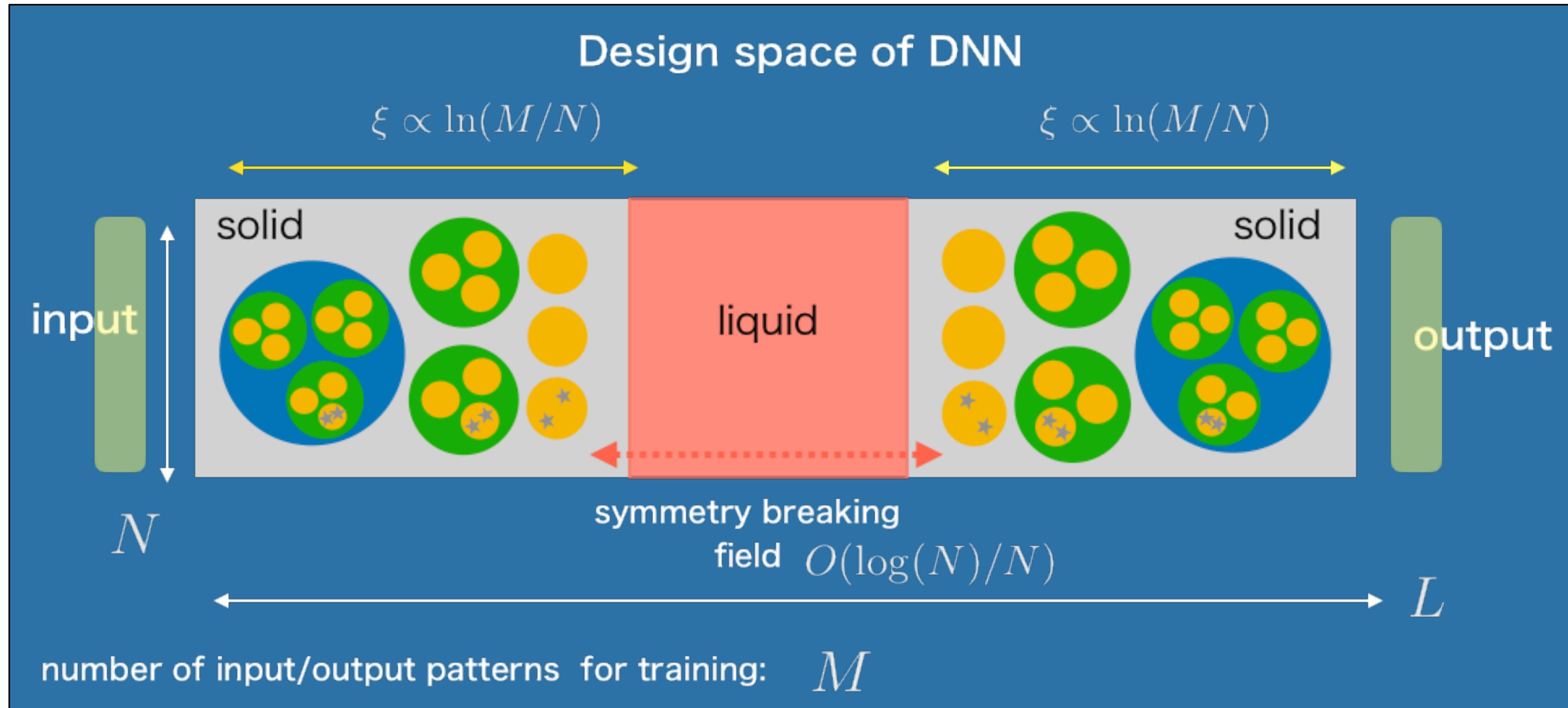
d)

$Q(x)$



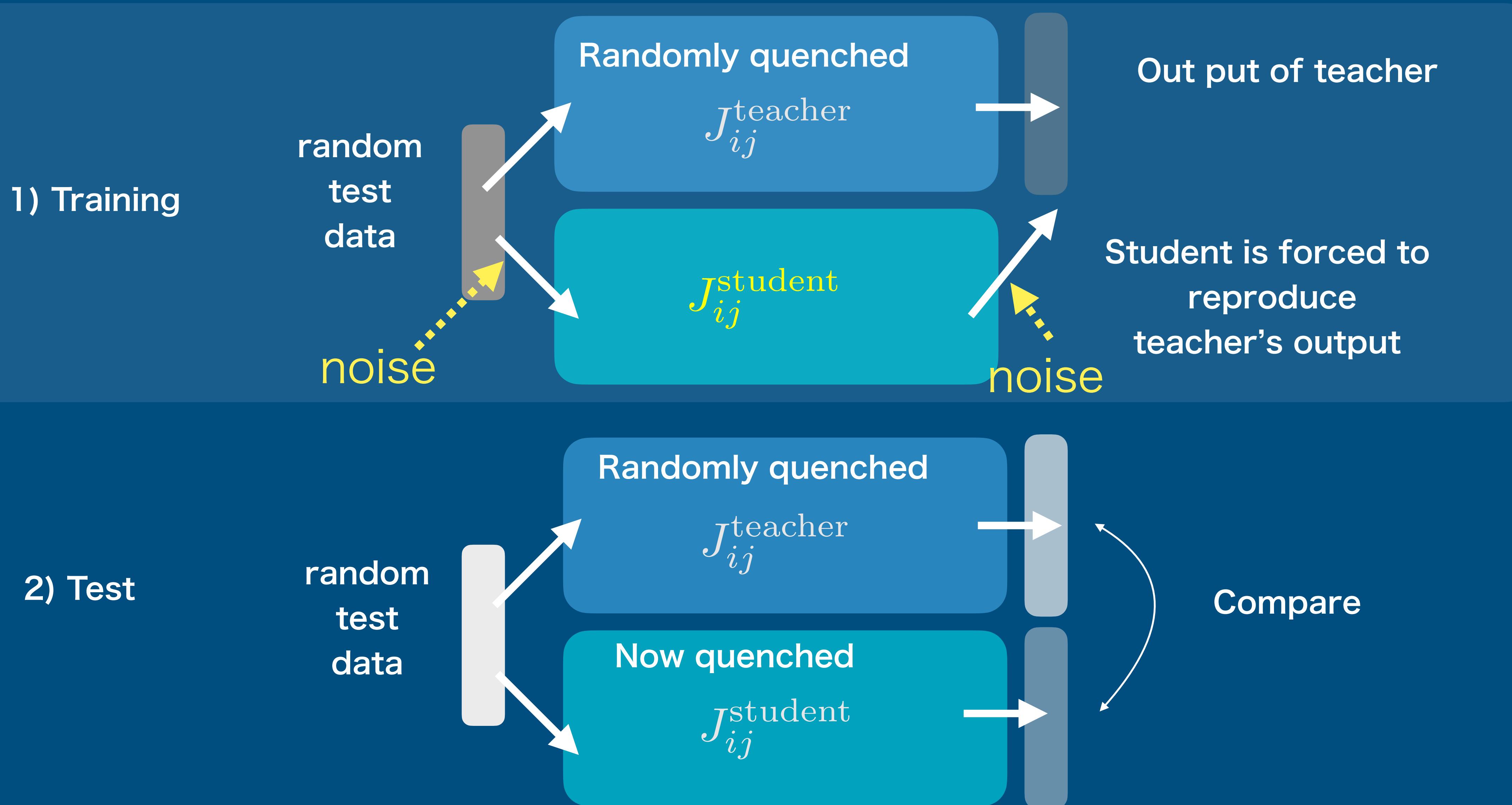
f)





each DNN is NOT a glass

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



Replicated Gardner volume

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

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)

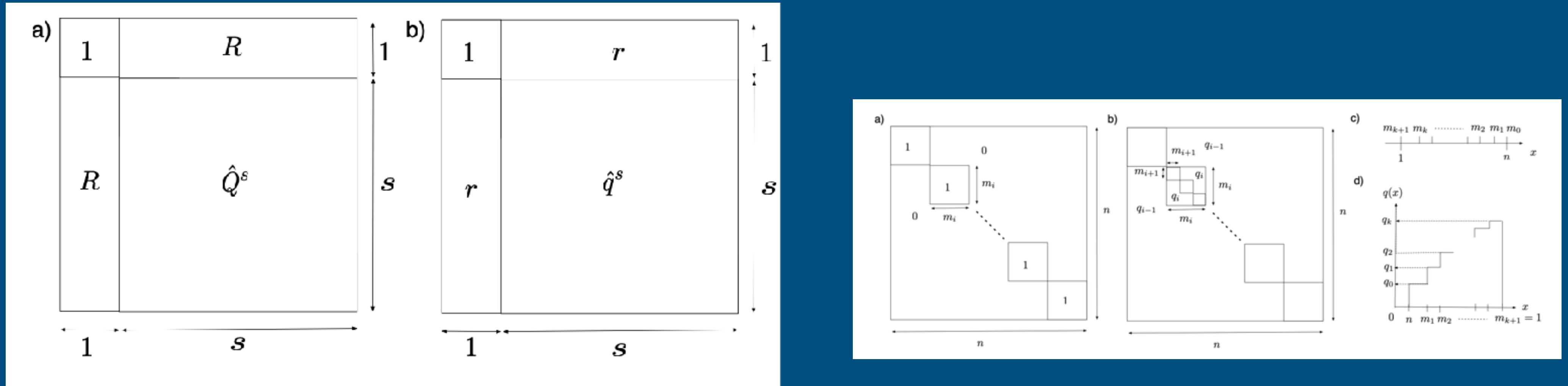
$$\frac{-\beta \overline{F(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} = \frac{\partial_n \overline{V^{\textcolor{blue}{1+n}}(\mathbf{S}_0, \mathbf{S}_L)}^{\text{visible}}}{NM} \Big|_{n=0} = \partial_n S_{\textcolor{blue}{1}+n}[\{\hat{Q}(l), \hat{q}(l)\}] \Big|_{n=0}$$

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

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

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

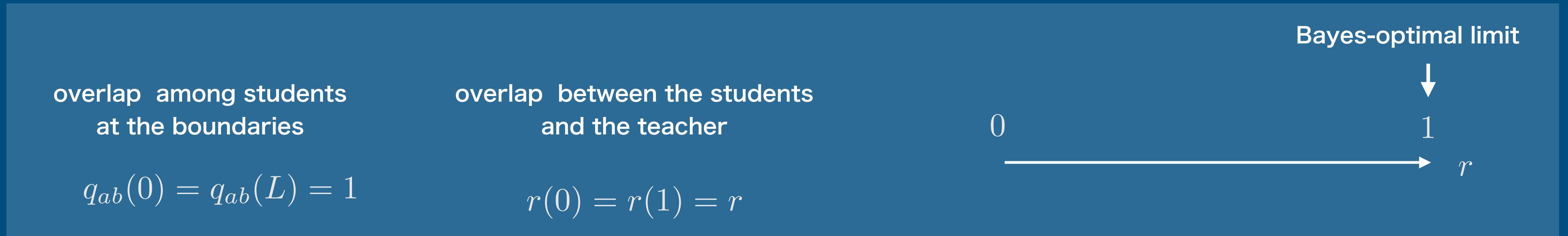
■ Parisi's RSB ansatz



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

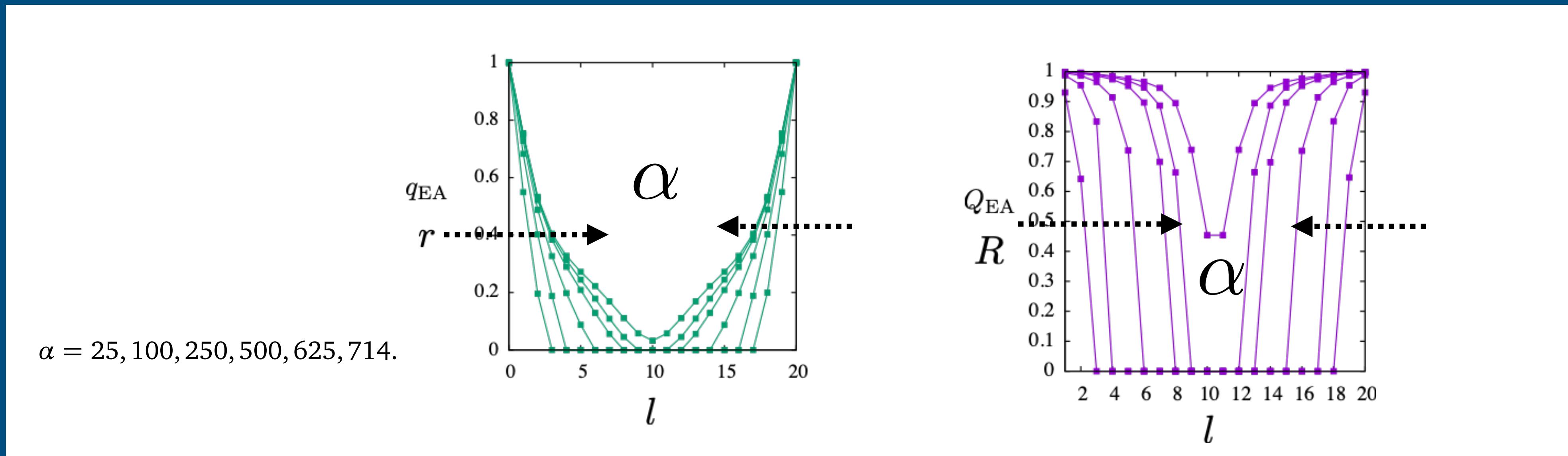
$$q_{ab}(l) = \sum_{i=0} q_i(l) (I_{ab}^{m_i} - I_{ab}^{m_{i+1}}) \quad l = 1, 2, \dots, L-1$$

■ Input/output boundaries



Bayes-optimal case (no noise)

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



“Wetting transition”

layer-by-layer 2nd order transition

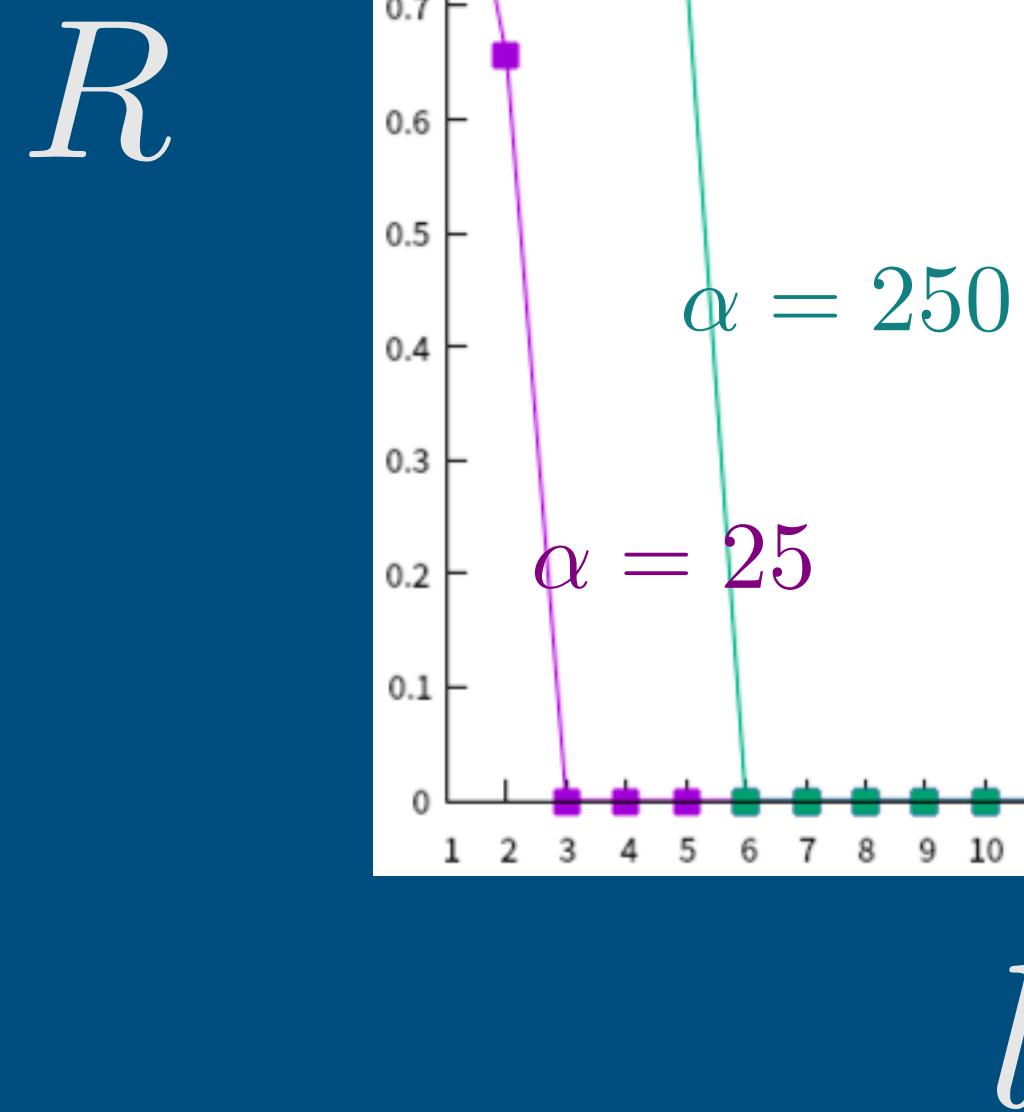
Bayes optimal, Nishimori condition

Replica symmetry (RS) holds

$q = r, Q = R$

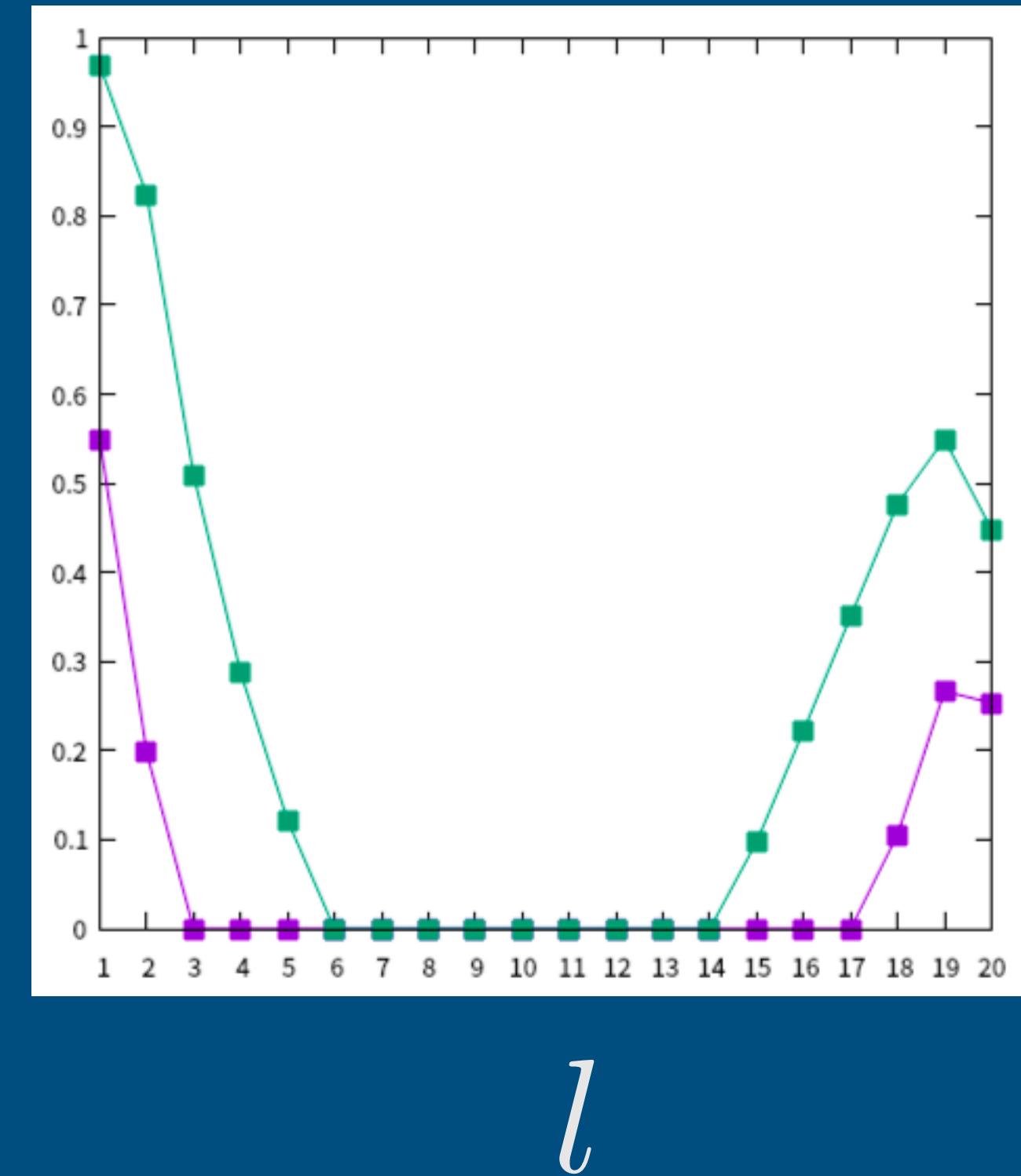
Generalization ability

teacher-student overlap of bond
after training



r_{test}

teacher-student overlap
for “test data”



over-parametrized DNN with “central liquid region” generalizes!
(if symmetry breaking field helps).

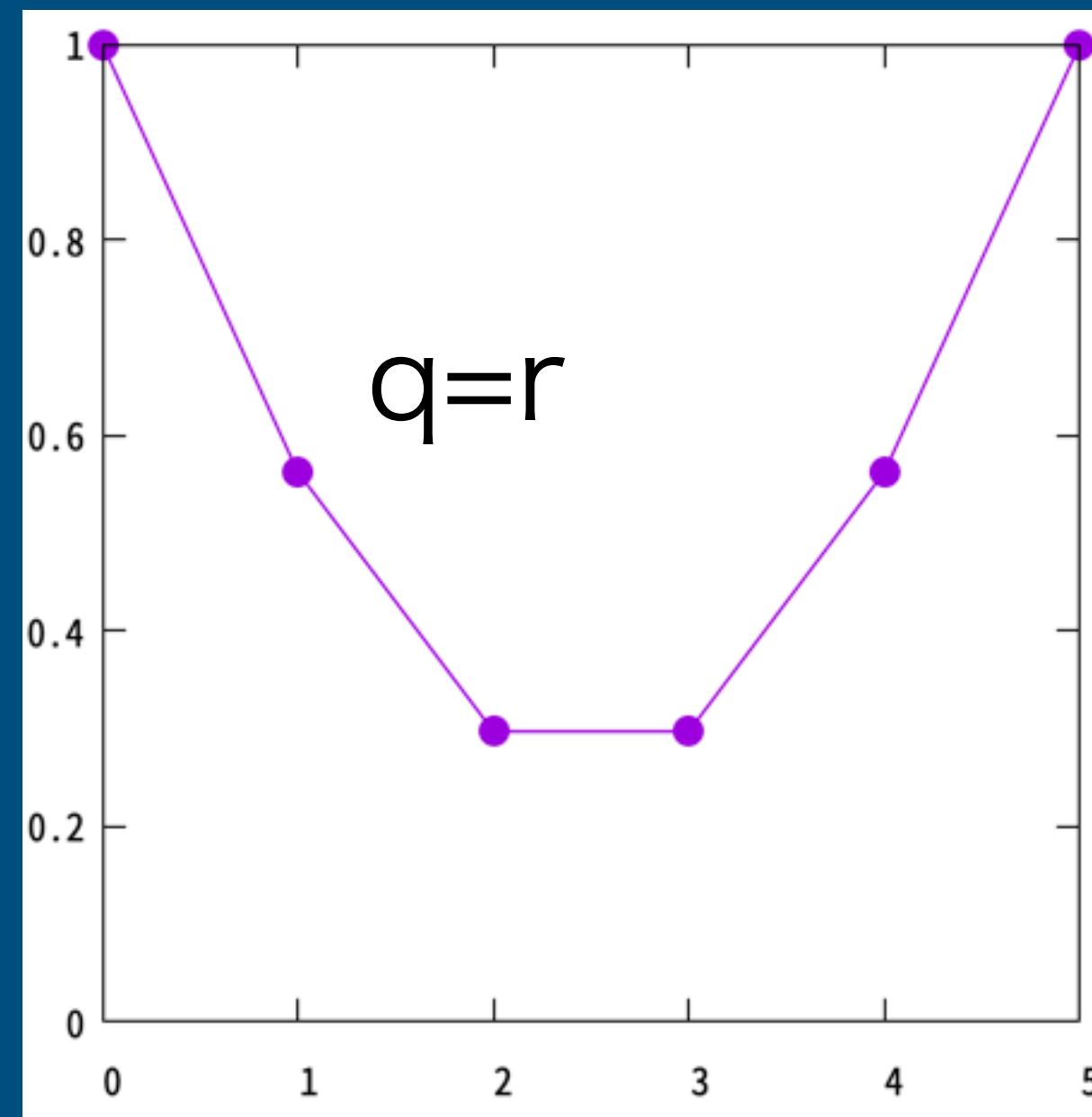
$r(l)$ is not fixed
non-zero solution!
on this side zero-
solution also exist

With noise

Spatial profile of the EA order parameter

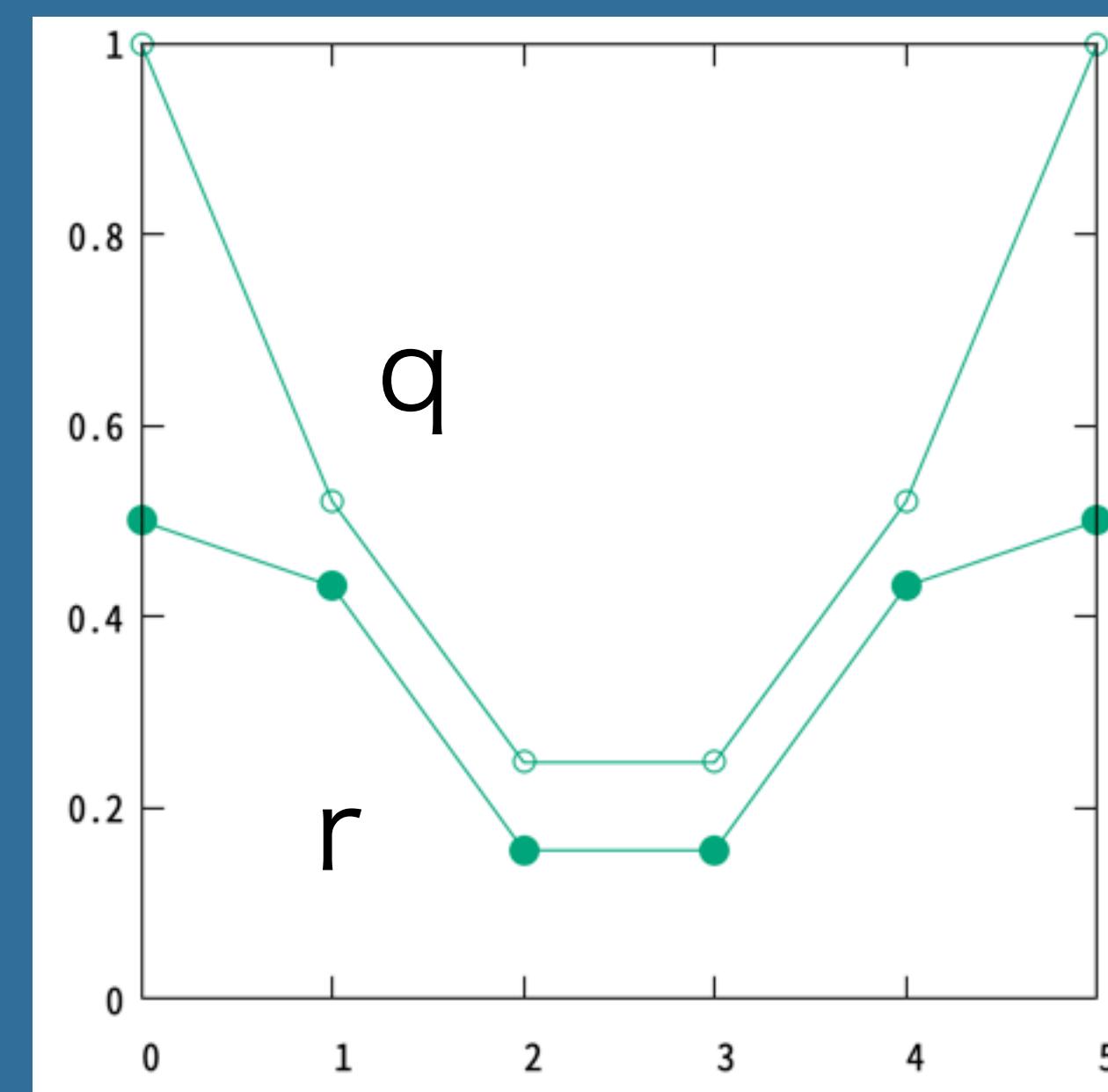
$$q_{\text{EA}} = q(1)$$

$$r = 1$$



l

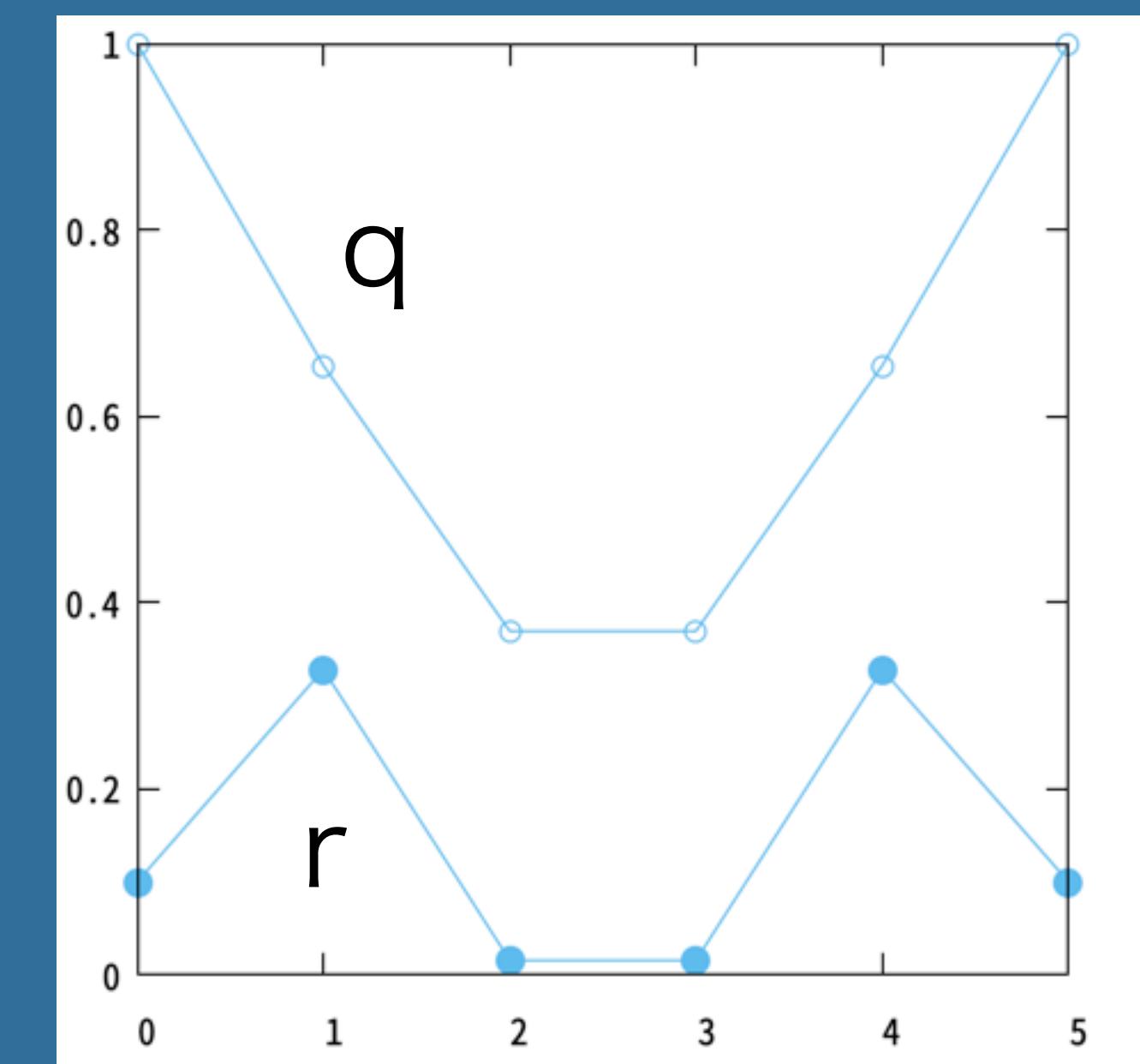
$$r = 0.5$$



l

noise

$$r = 0.1$$



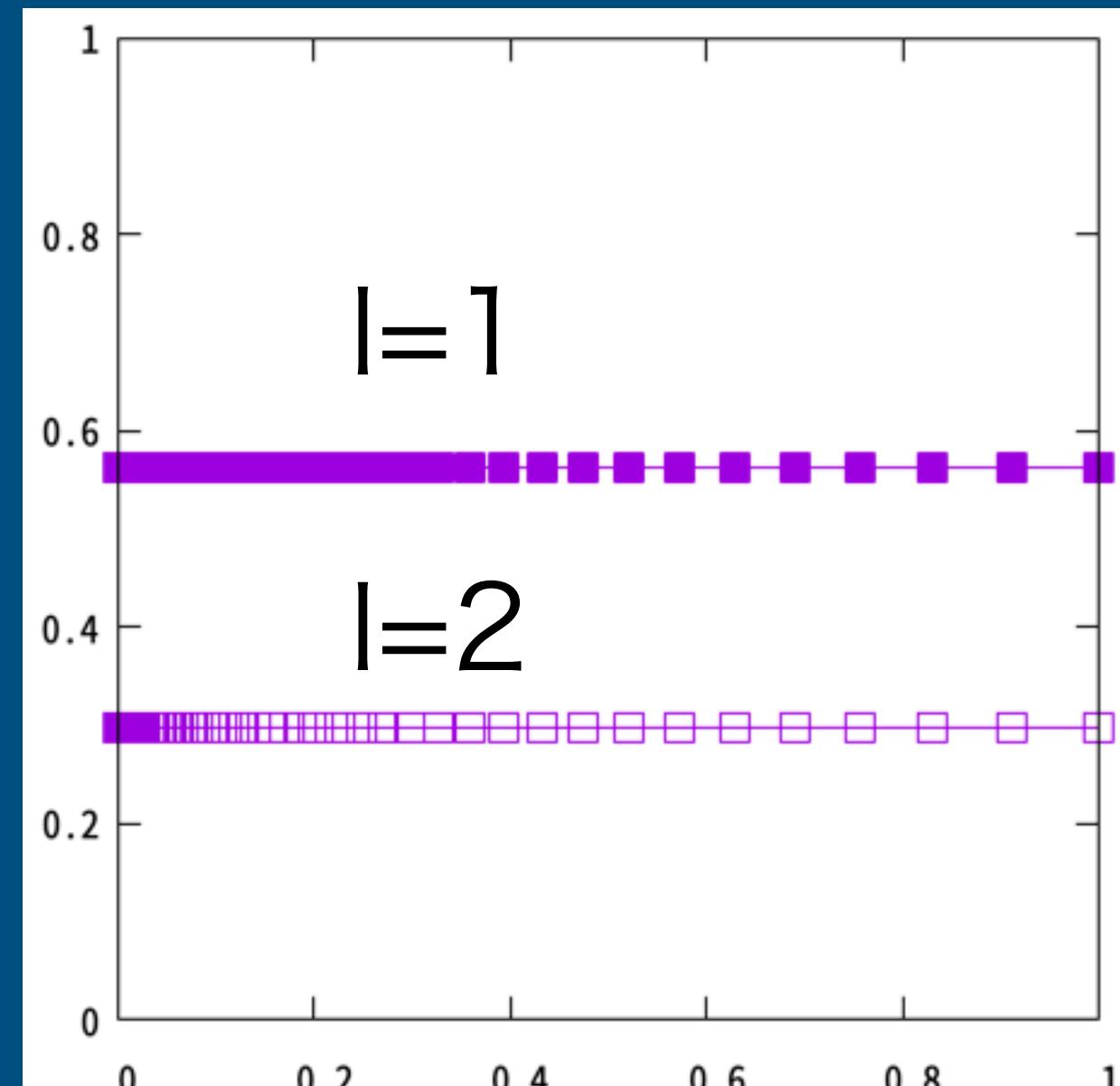
l

$$\alpha = 25$$

Hierarchical structure of the solutions

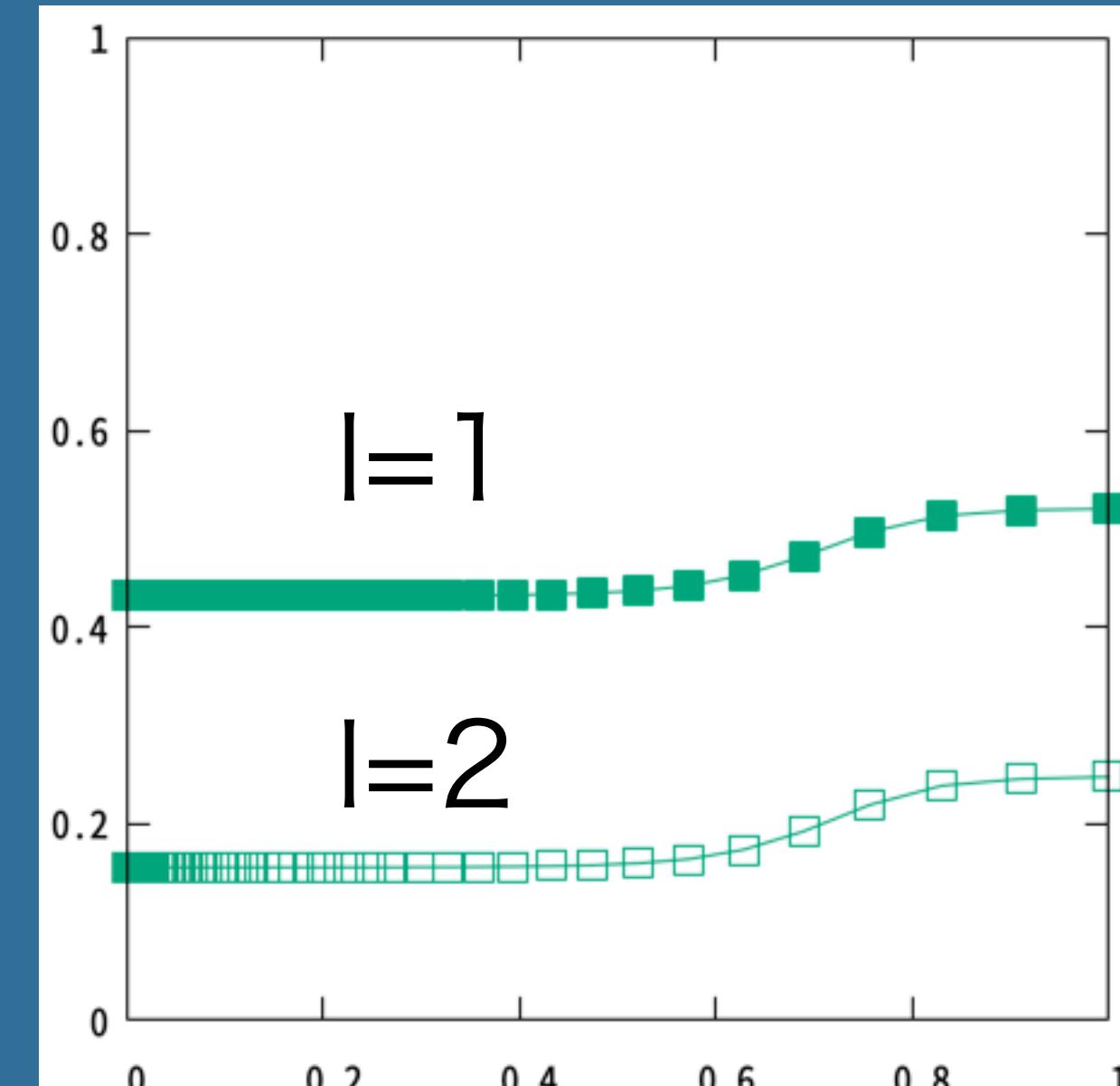
$$q(x)$$

$$r = 1$$

 x

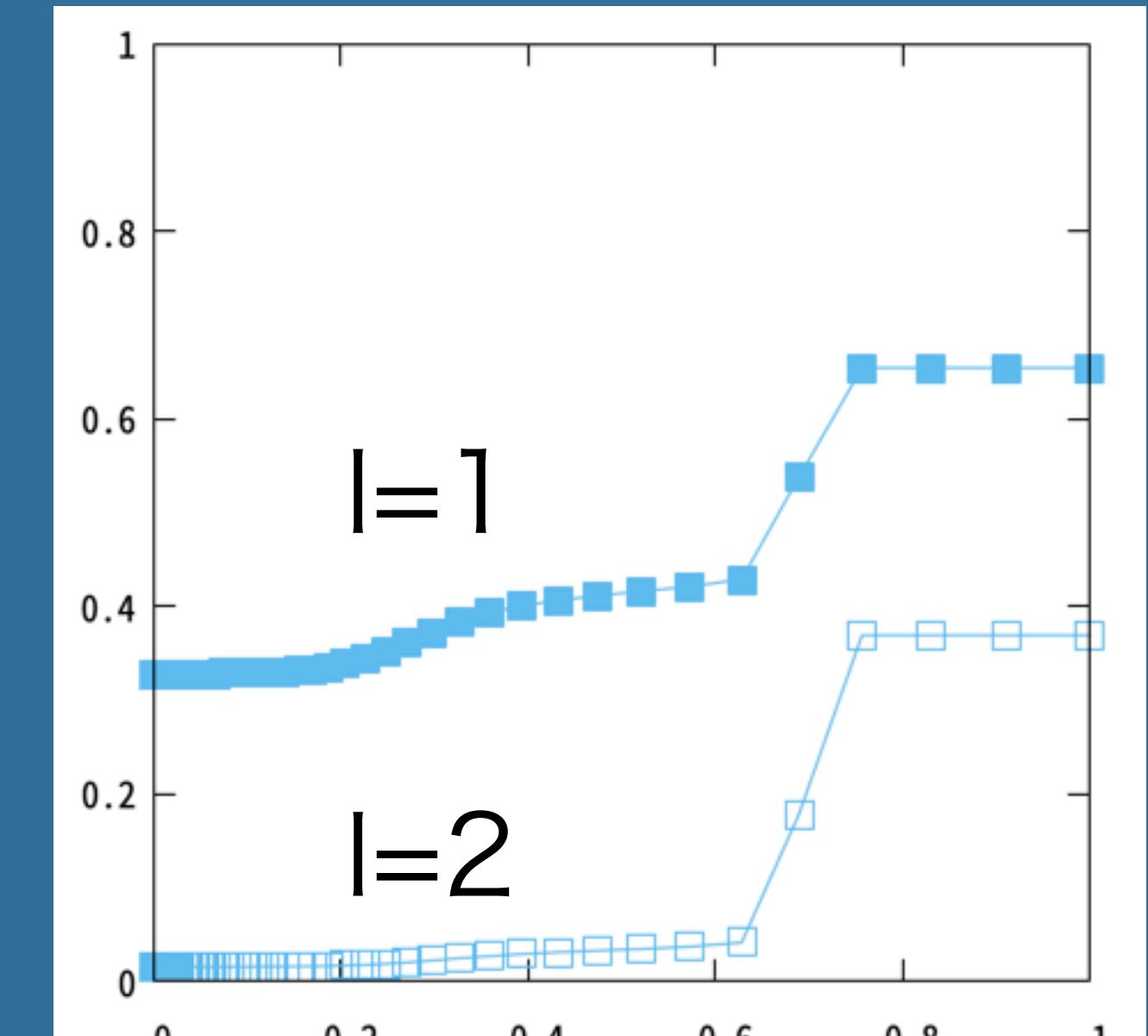
Replica symmetric

$$r = 0.5$$

 x

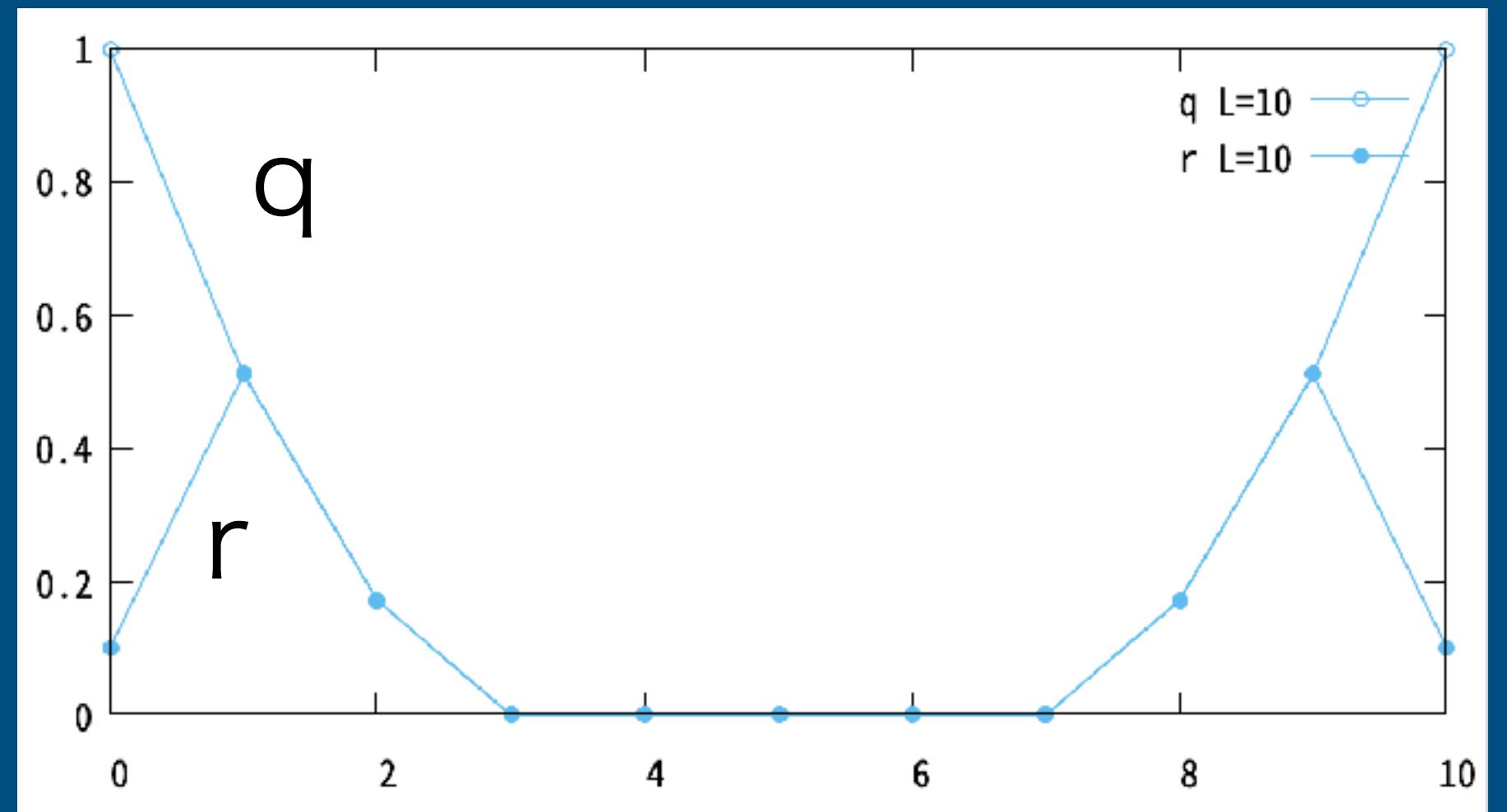
Replica symmetry broken

$$r = 0.1$$

 x

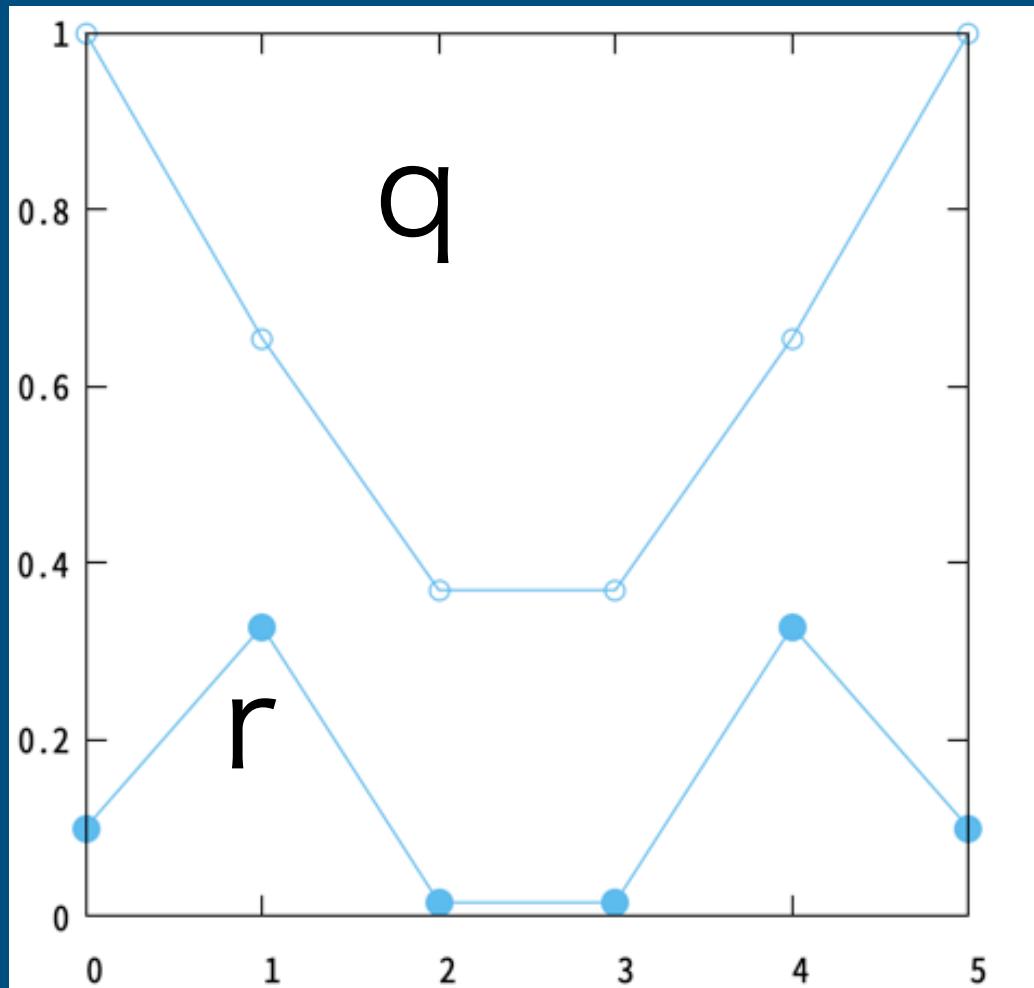
$r = 0.1$

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



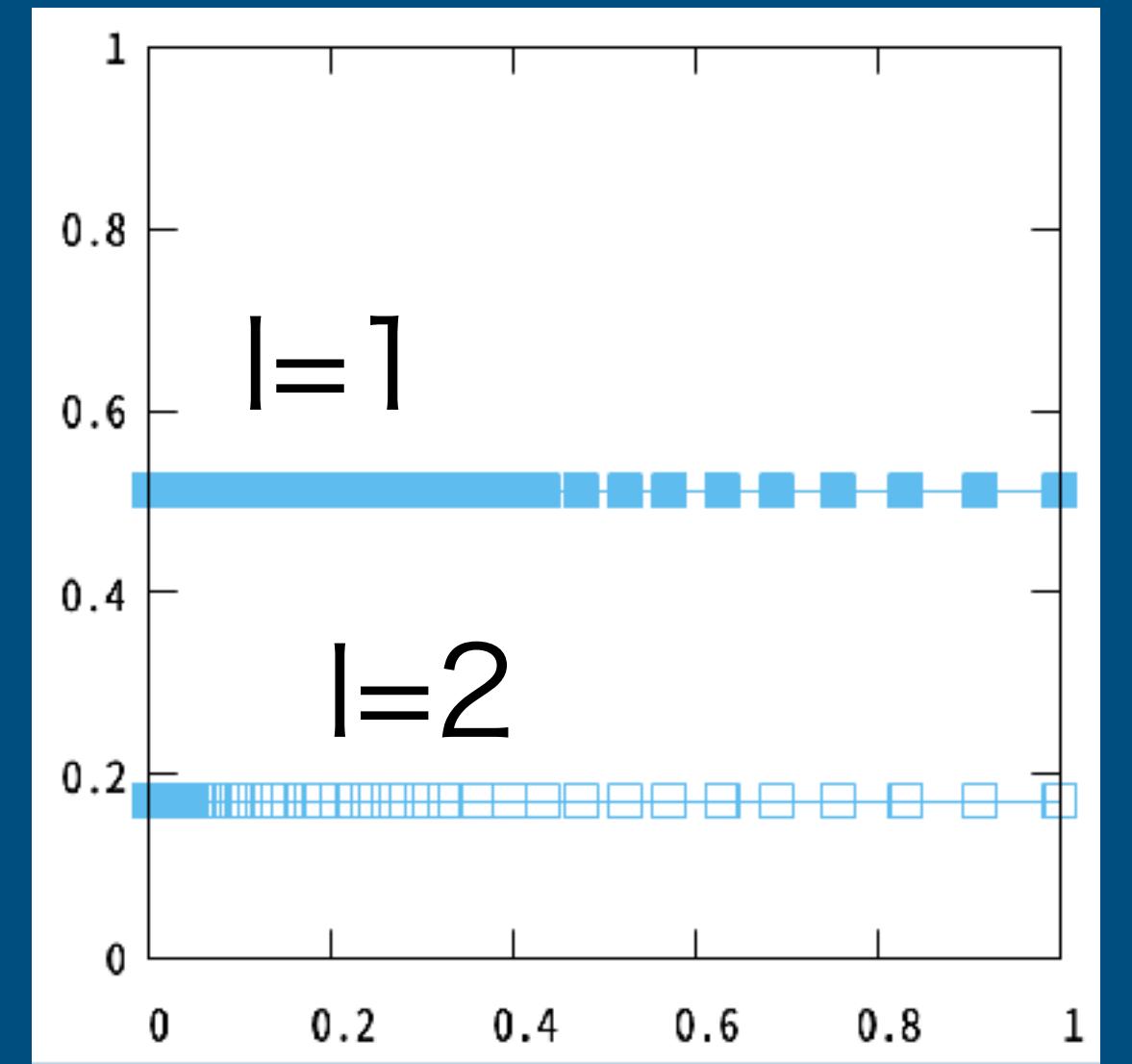
$r = 0.1$

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



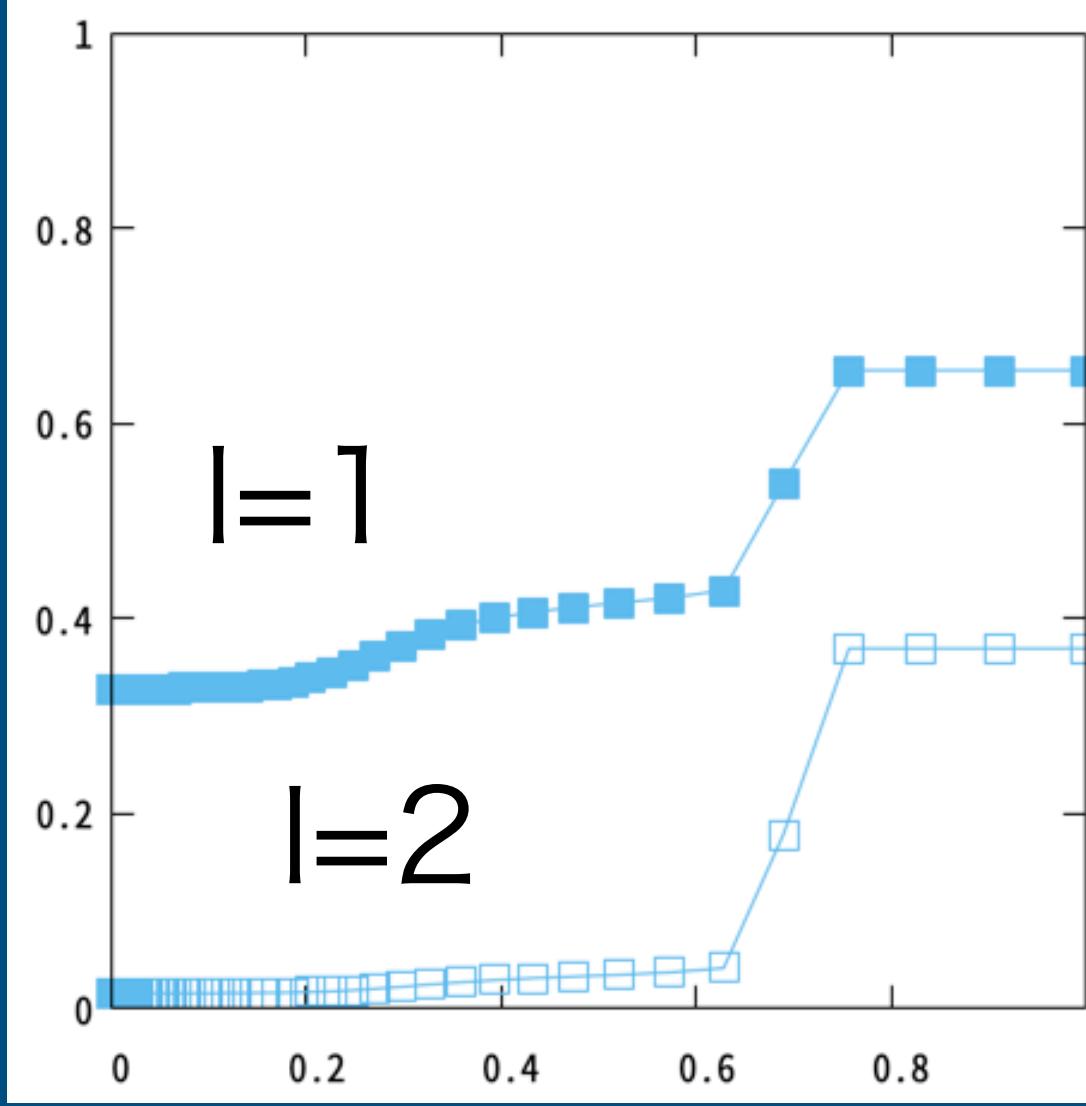
Replica symmetric!

$$q(x)$$



$$x$$

$$q(x)$$



$$x$$

■ Simulation of learning in a teacher-student setting

t

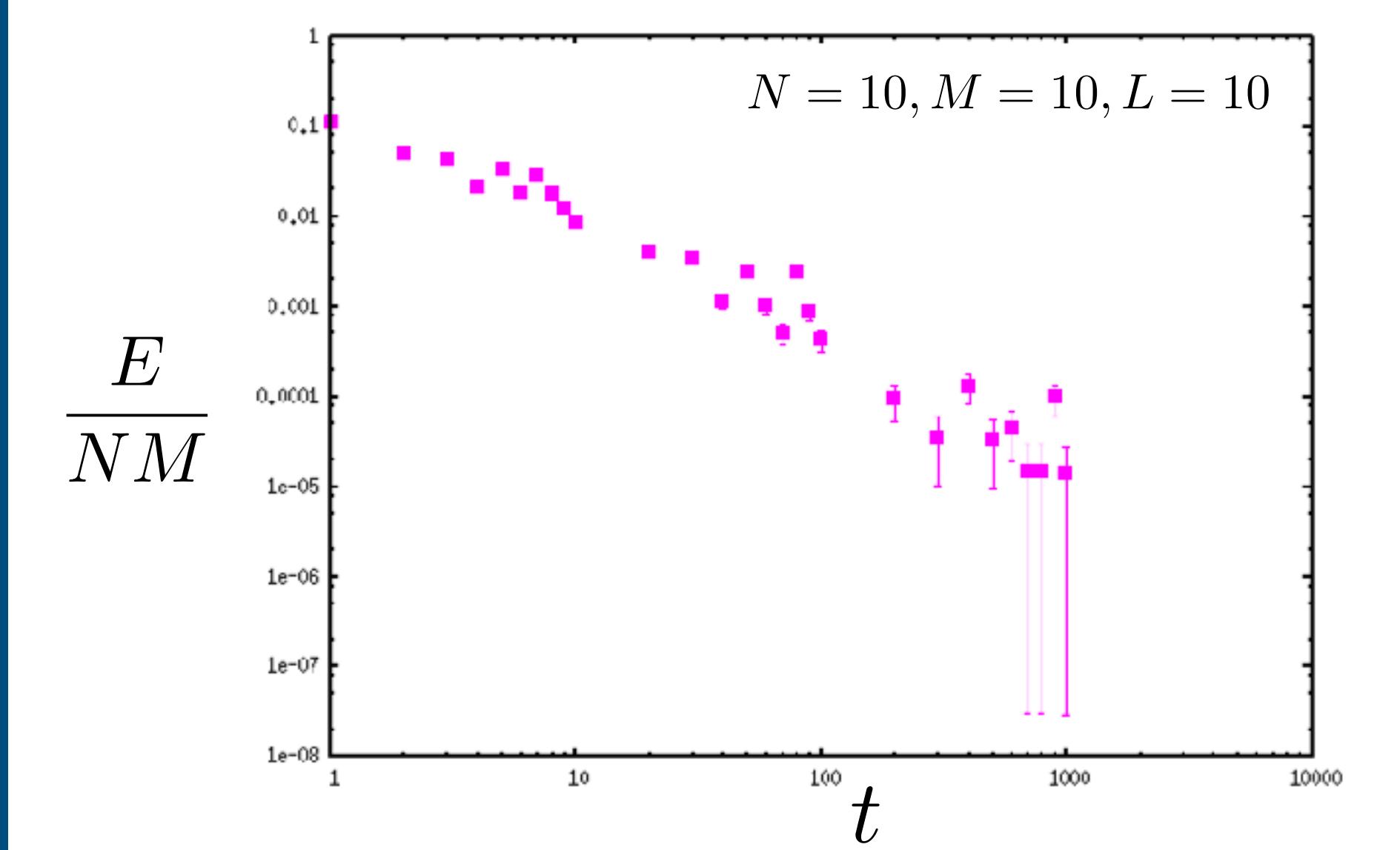
T=0 greedy Monte Carlo

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

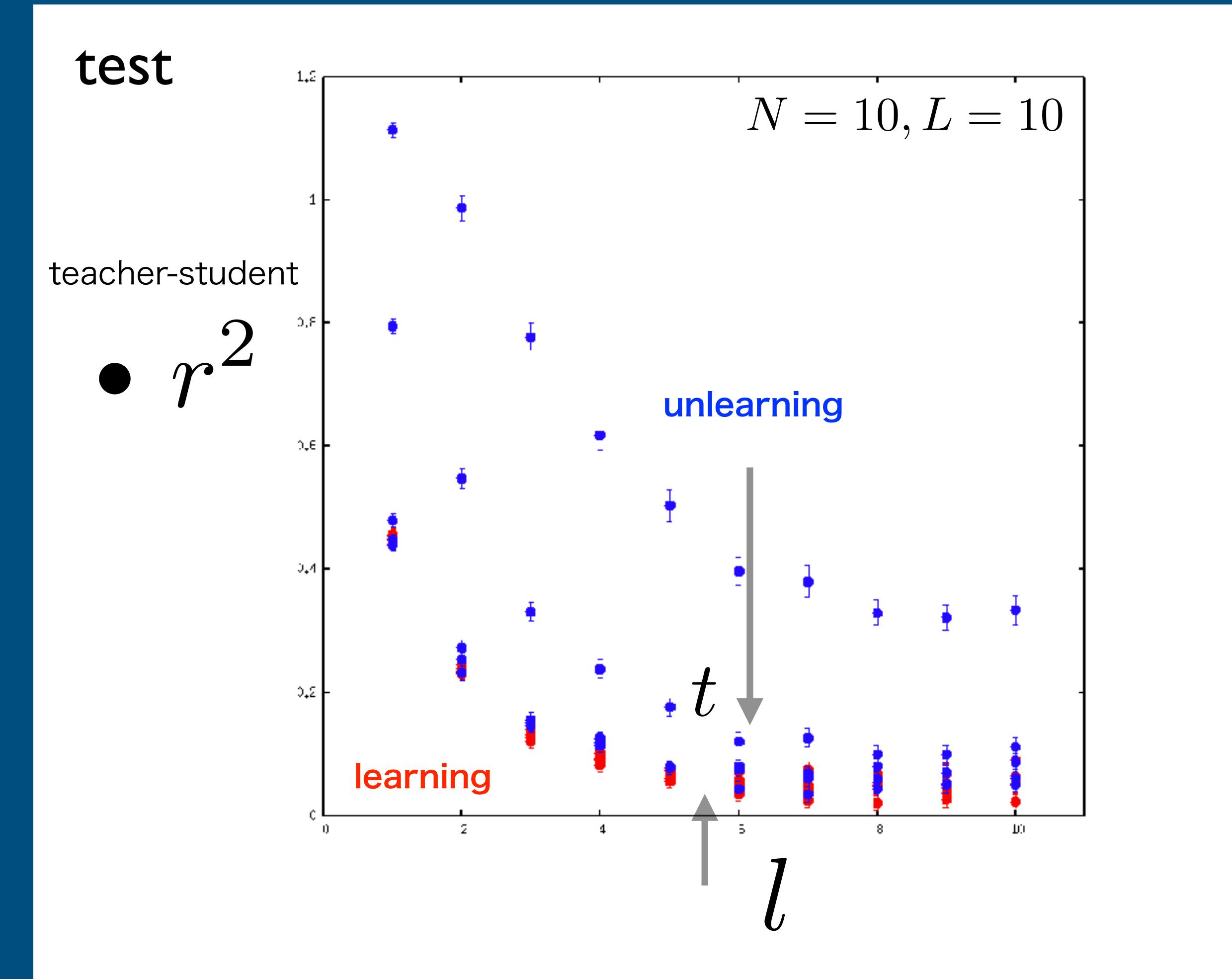
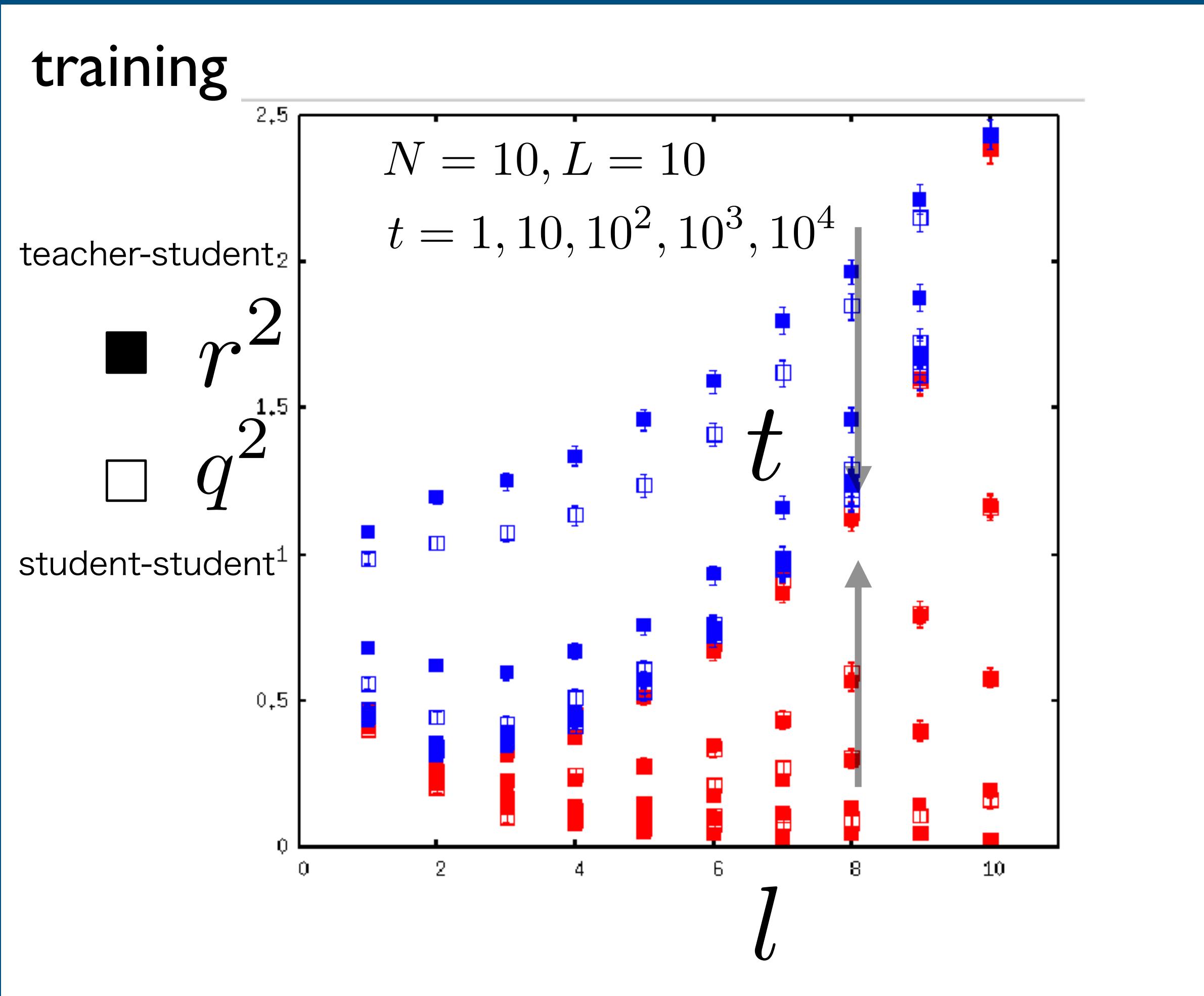
2. **“Learning”** : start from random configuration



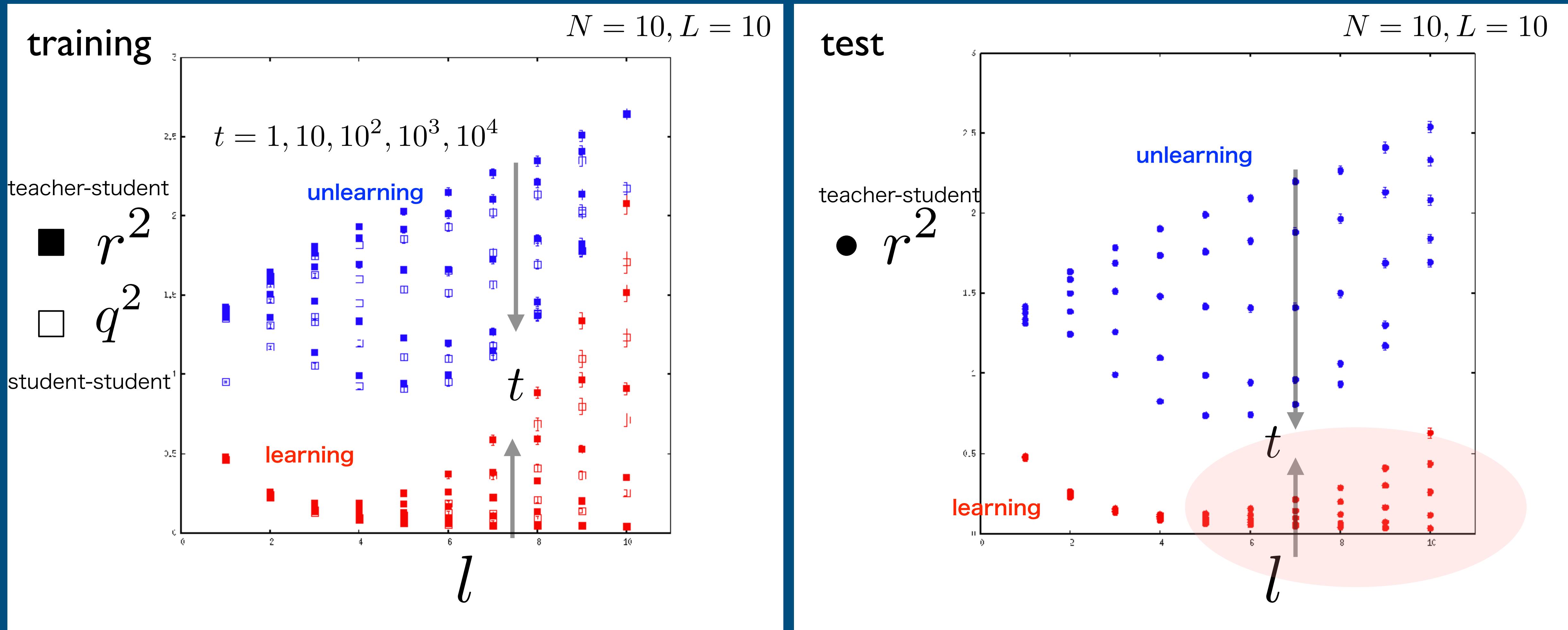
Permutation-invariant overlap

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

$$\alpha = 1$$



$$\alpha = 32$$



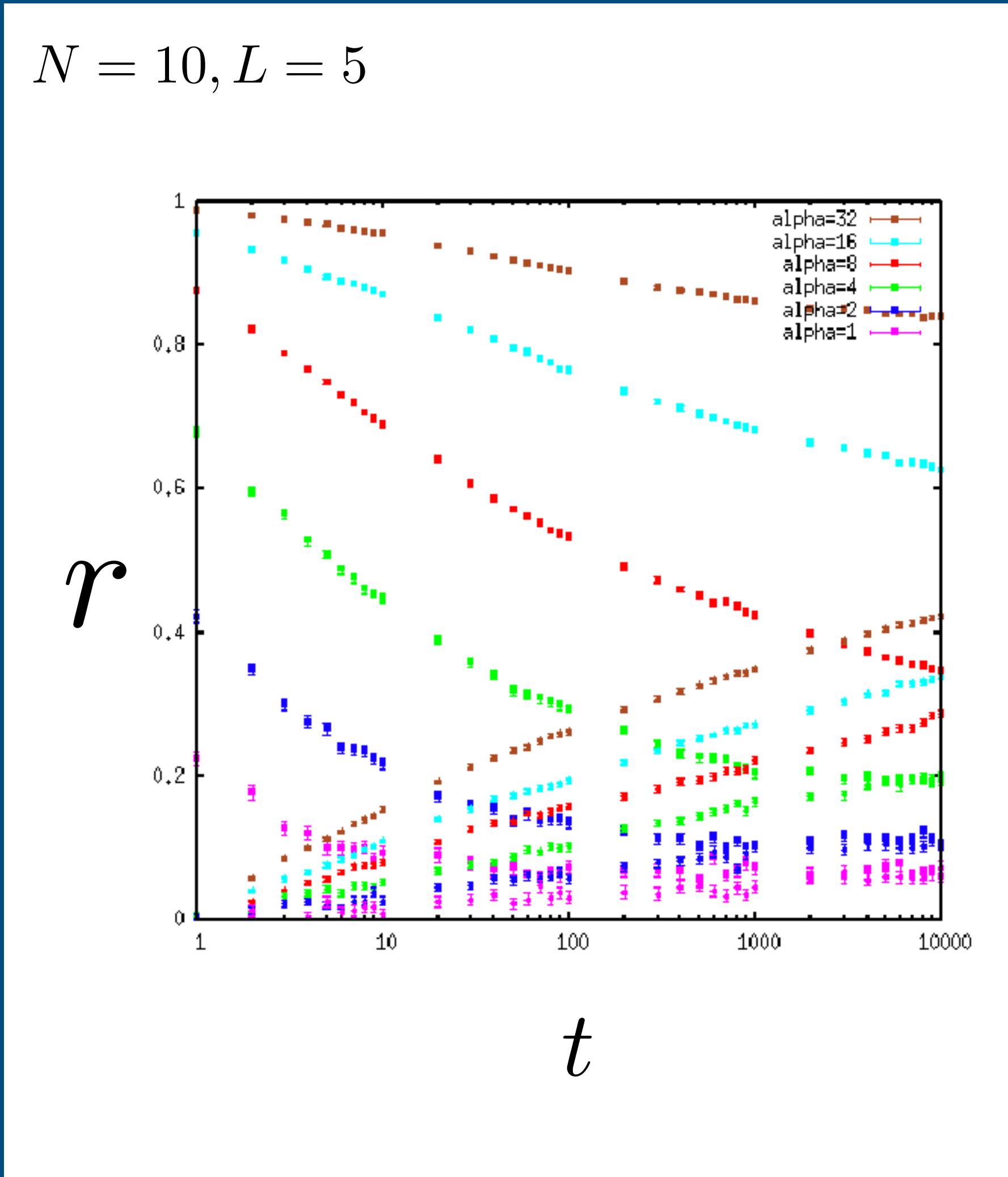
Generalization ability

unlearning

learning

teacher-student overlap in the test

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



Generalization ability

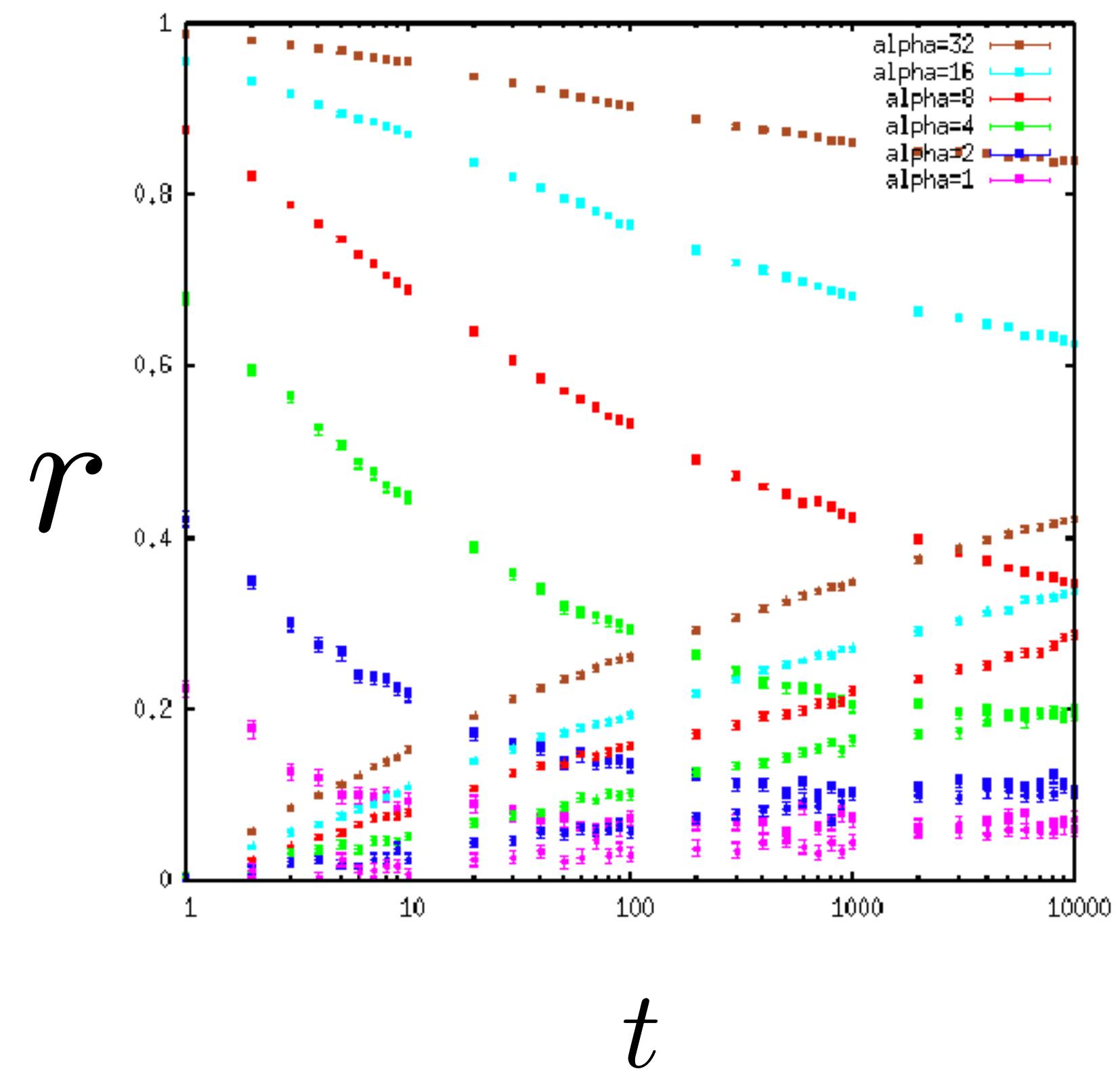
teacher-student overlap in the test

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

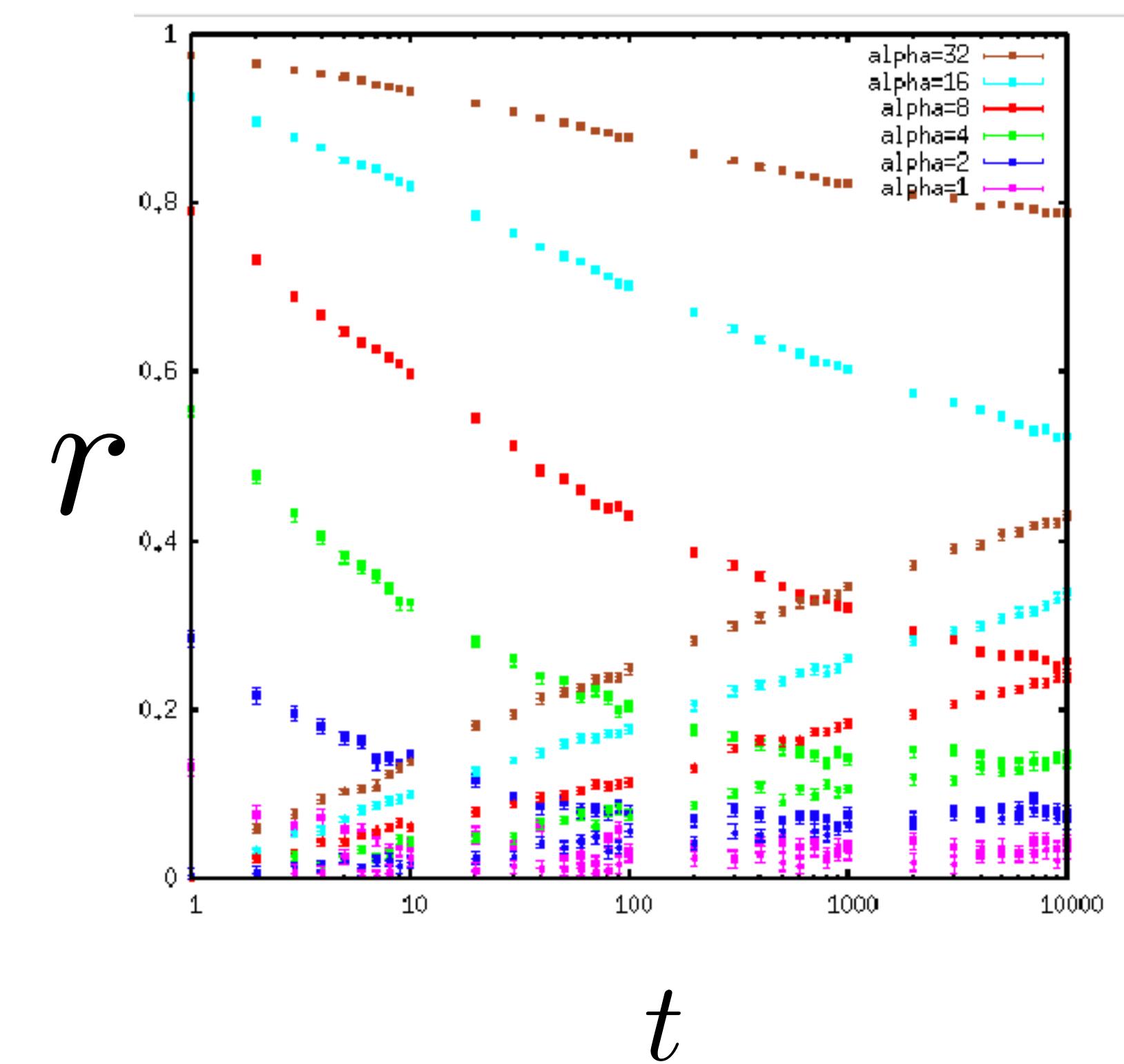
unlearning

learning

$N = 10, L = 5$



$N = 10, L = 10$



deeper systems generalize as well...

■ 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: loop corrections...

mismatch of architecture

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

other activation functions: sigmoid, ReLU,....

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