

Machine learning: from mammal's brain to statistical mechanics

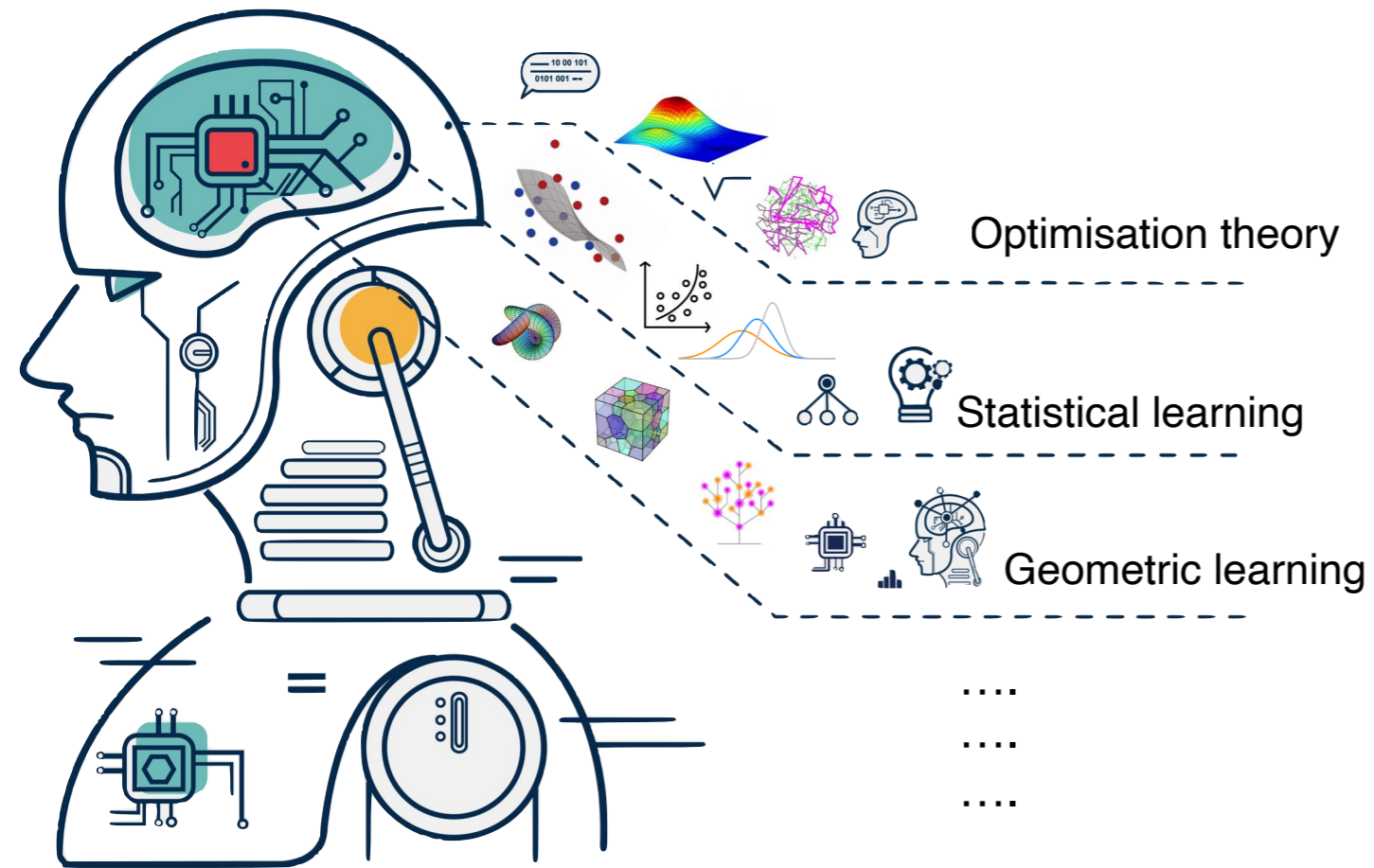
Elena Agliari

Dipartimento di Matematica, Sapienza Università di Roma
Istituto Nazionale di Alta Matematica, Roma



Spring School “Bruno Touschek”
Laboratori Nazionali Frascati
16 May 2024

Hard sciences for machine learning

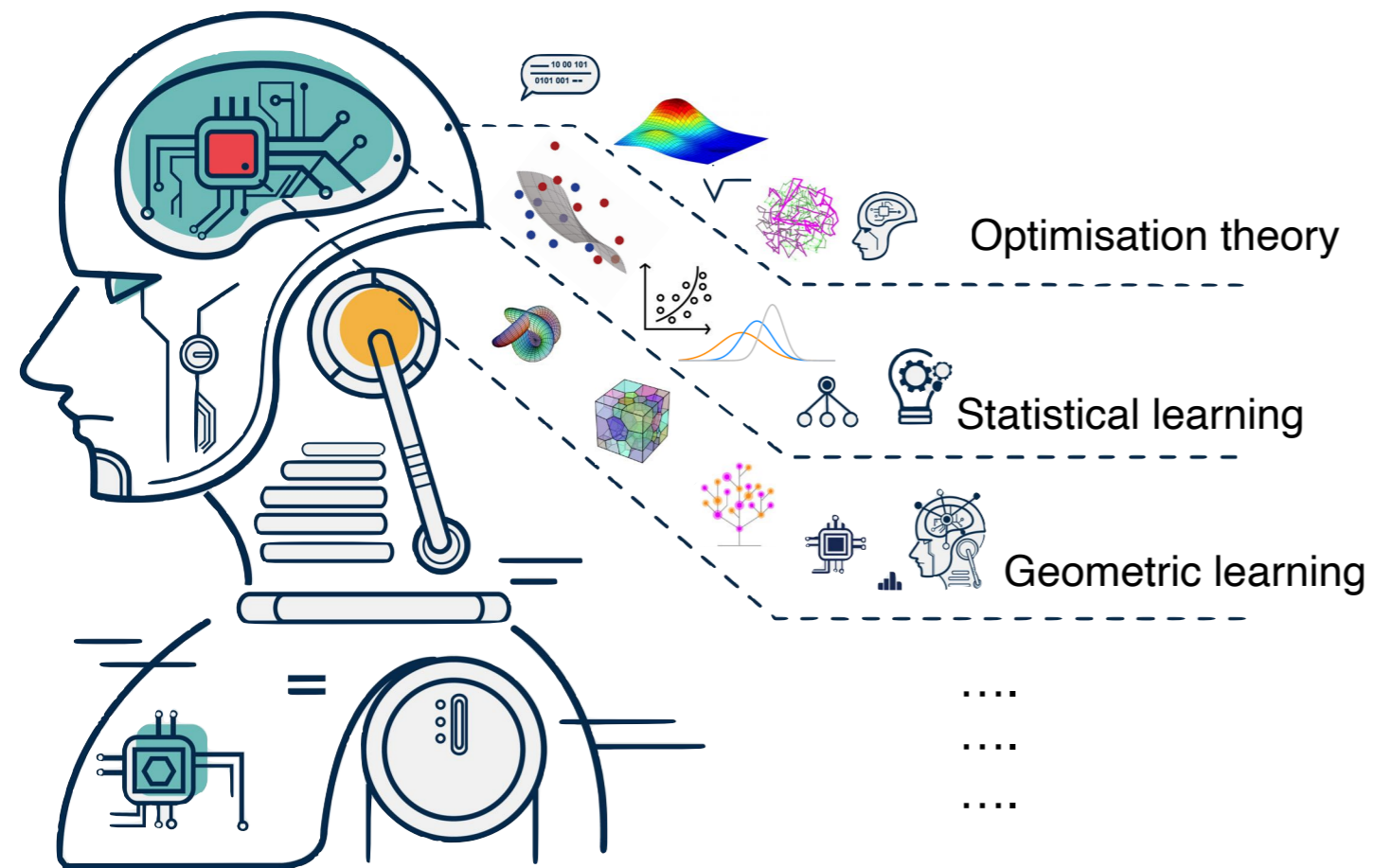


Mathematical control

- design optimal architecture
- estimate hyperparameters
- check dataset size
-

→ Sustainability & Interpretability

Hard sciences for machine learning

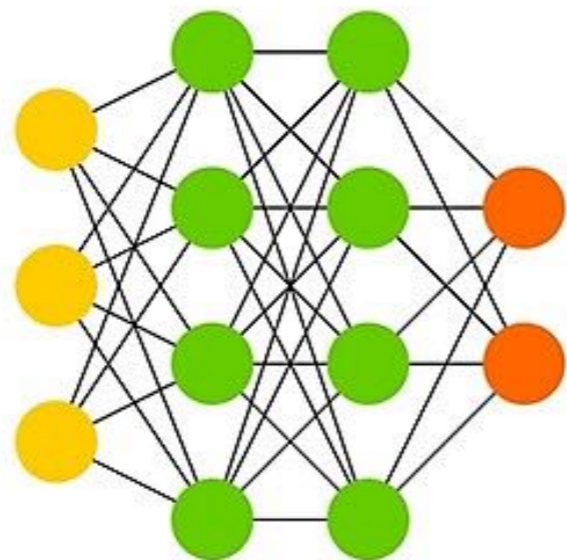


Mathematical control

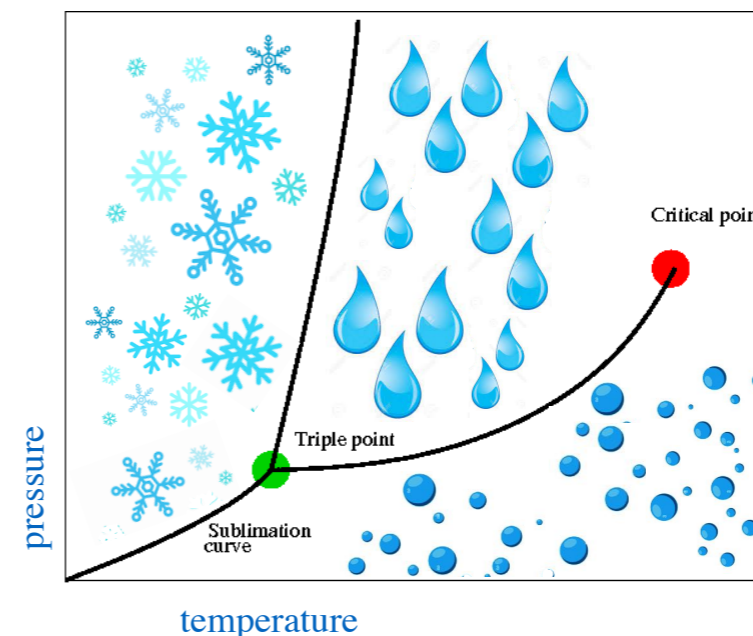
- design optimal architecture
- estimate hyperparameters
- check dataset size
-

→ Sustainability & Interpretability

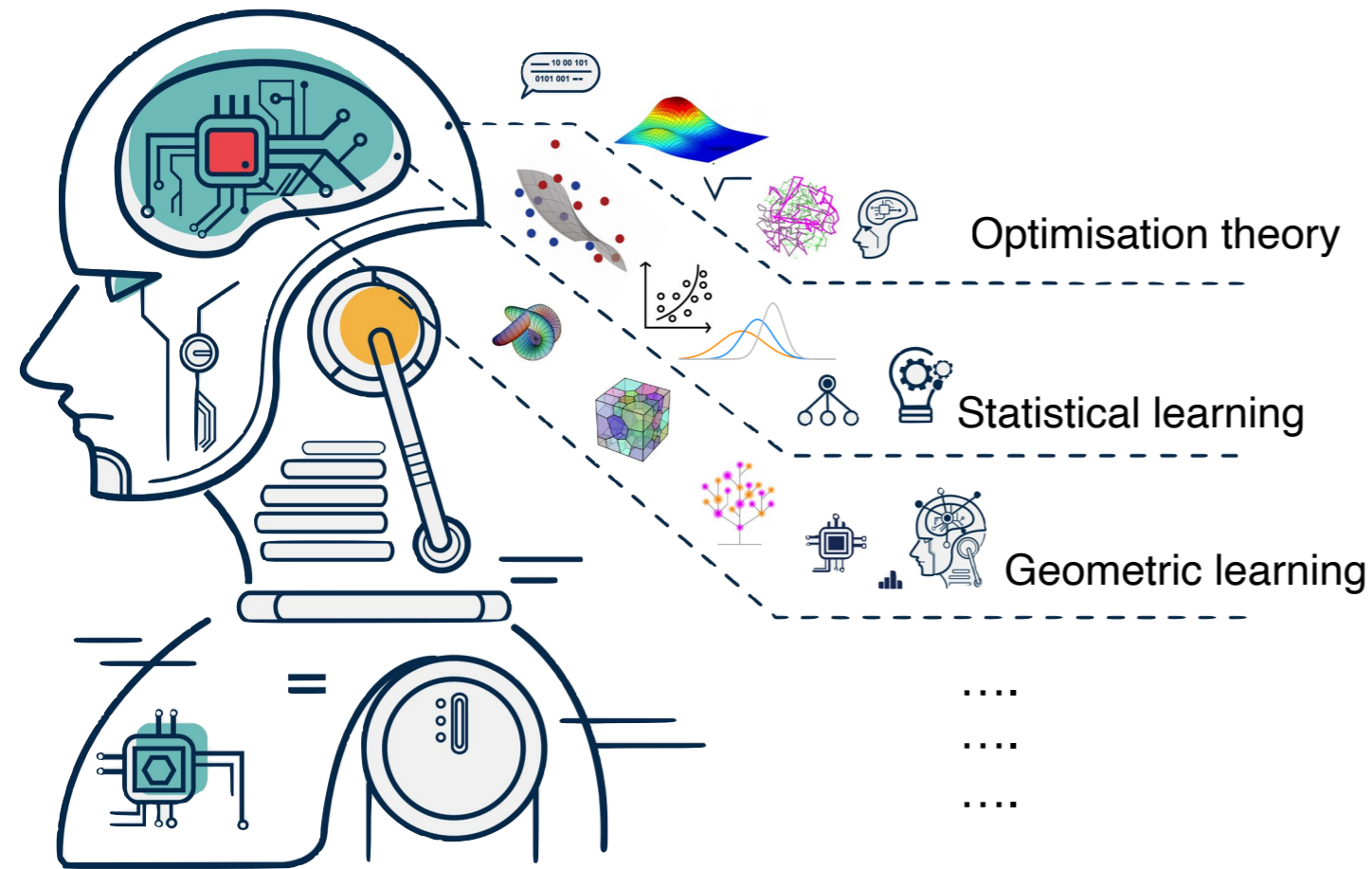
Statistical-mechanics perspective



Information-processing capability as emerging collective behavior



Hard sciences for machine learning

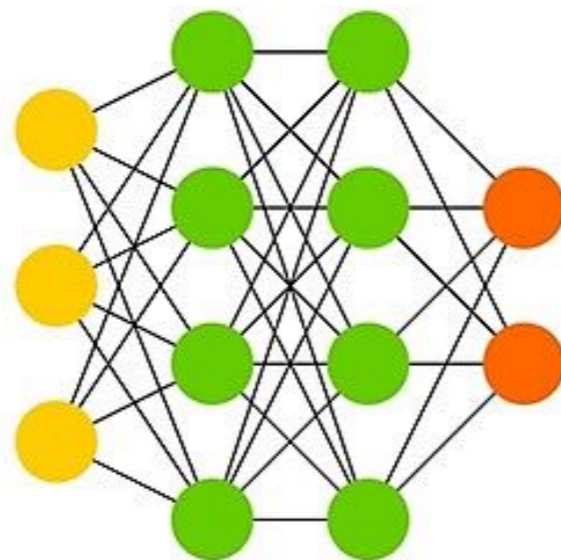


Mathematical control

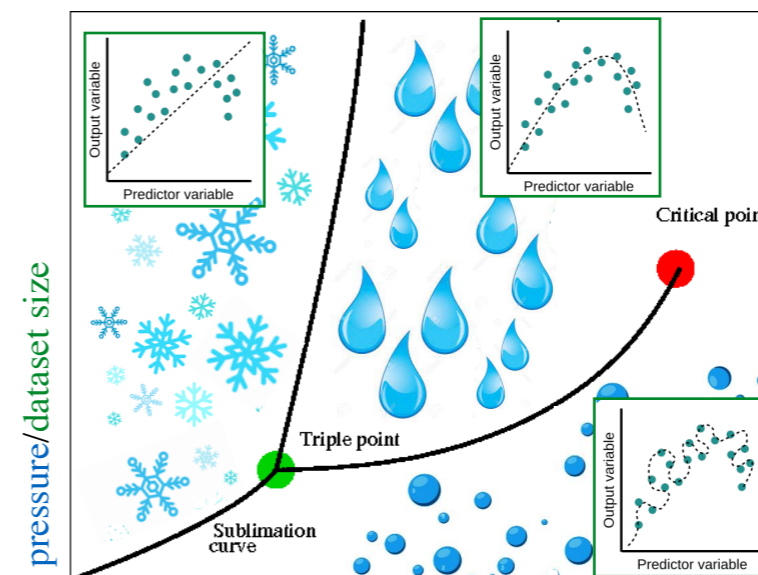
- design optimal architecture
- estimate hyperparameters
- check dataset size
-

→ Sustainability & Interpretability

Statistical-mechanics perspective



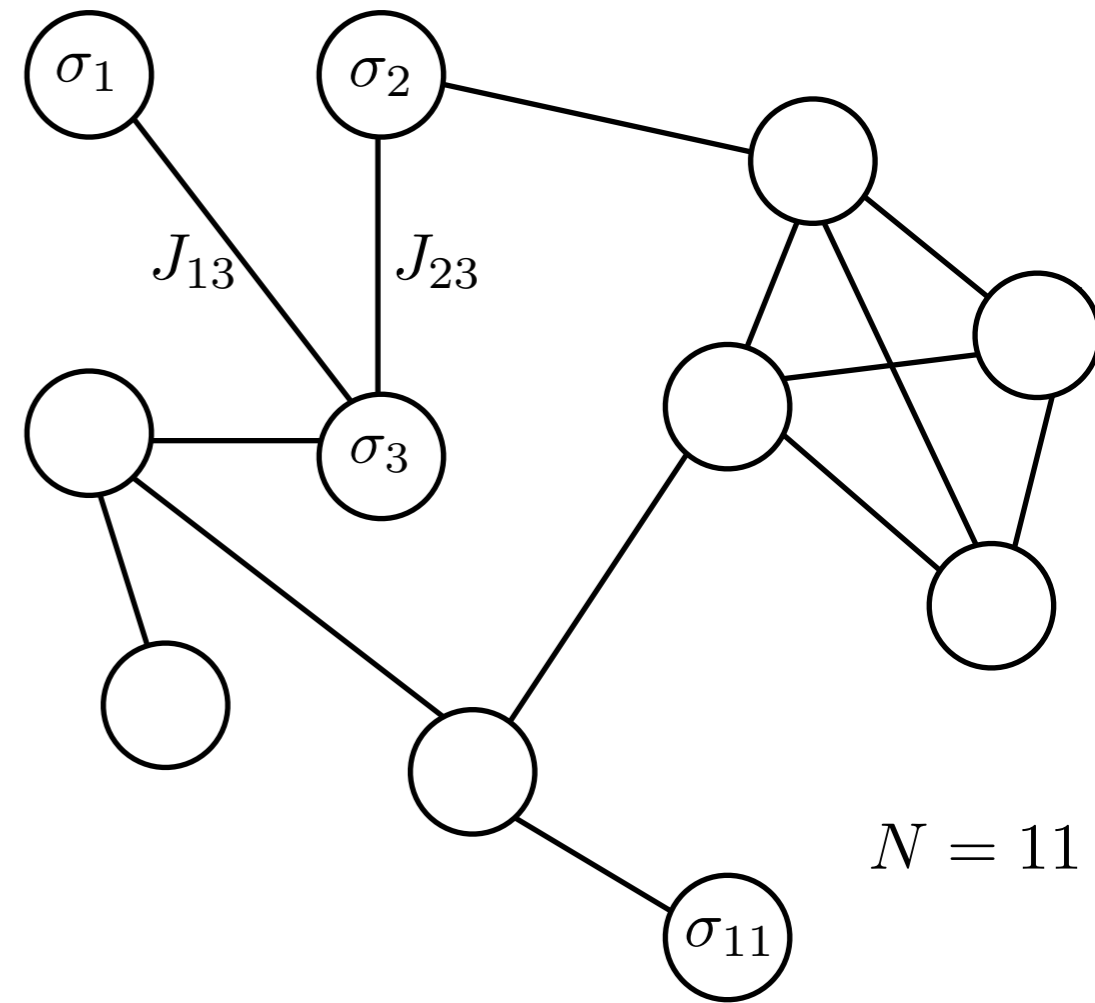
Information-processing capability as emerging collective behaviour



temperature/depth

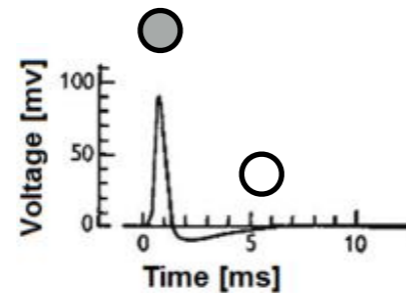
Neural networks: a biology-inspired introduction

Nodes as neurons with activity $\sigma = \{\sigma_i\}_{i=1,\dots,N}$
 Links as synapses with efficacy $\mathbf{J} = \{J_{ij}\}_{i,j=1,\dots,N}$



Firing/quiescent neurons

$$\sigma \in \{-1, +1\}^N$$

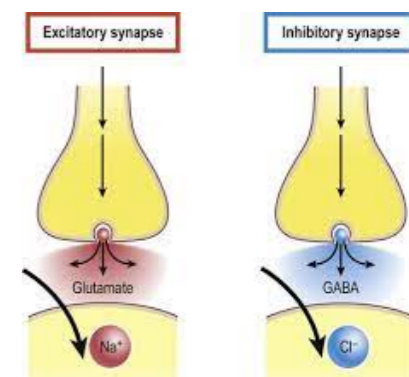


$\{\sigma_i\}$ fast degrees of freedom

Excitatory/Inhibitory synapses

$$\mathbf{J} \in \mathbb{R}^{N \times N}$$

$\{J_{ij}\}$ slow degrees of freedom

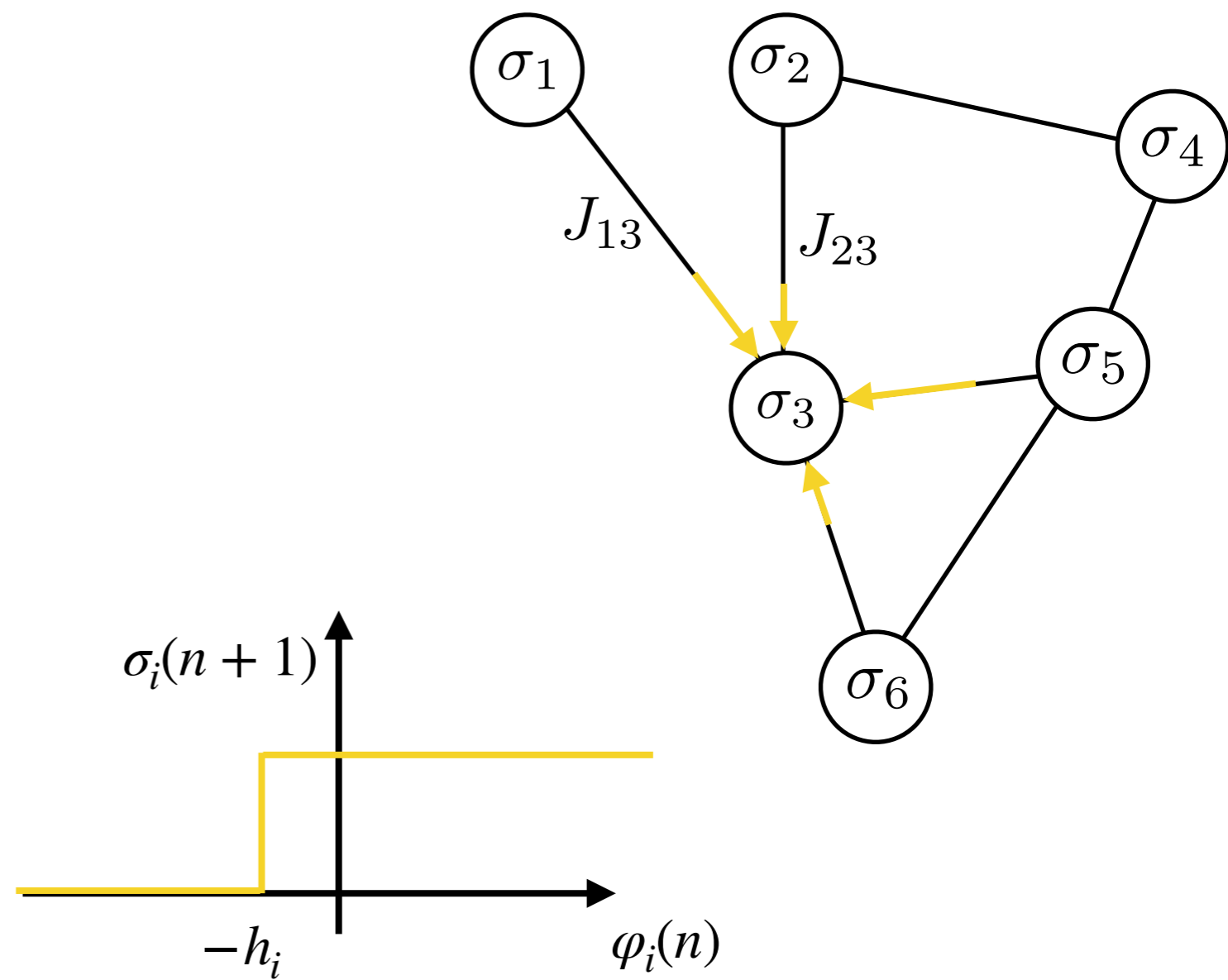


Neural dynamics

$$\sigma_i(n + 1) = \text{sgn}(\varphi_i(n) + h_i)$$

Field acting on i at time step n Firing threshold

$$\varphi_i(n) = \sum_{\substack{k=1 \\ k \neq i}}^N J_{ik} \sigma_k(n)$$

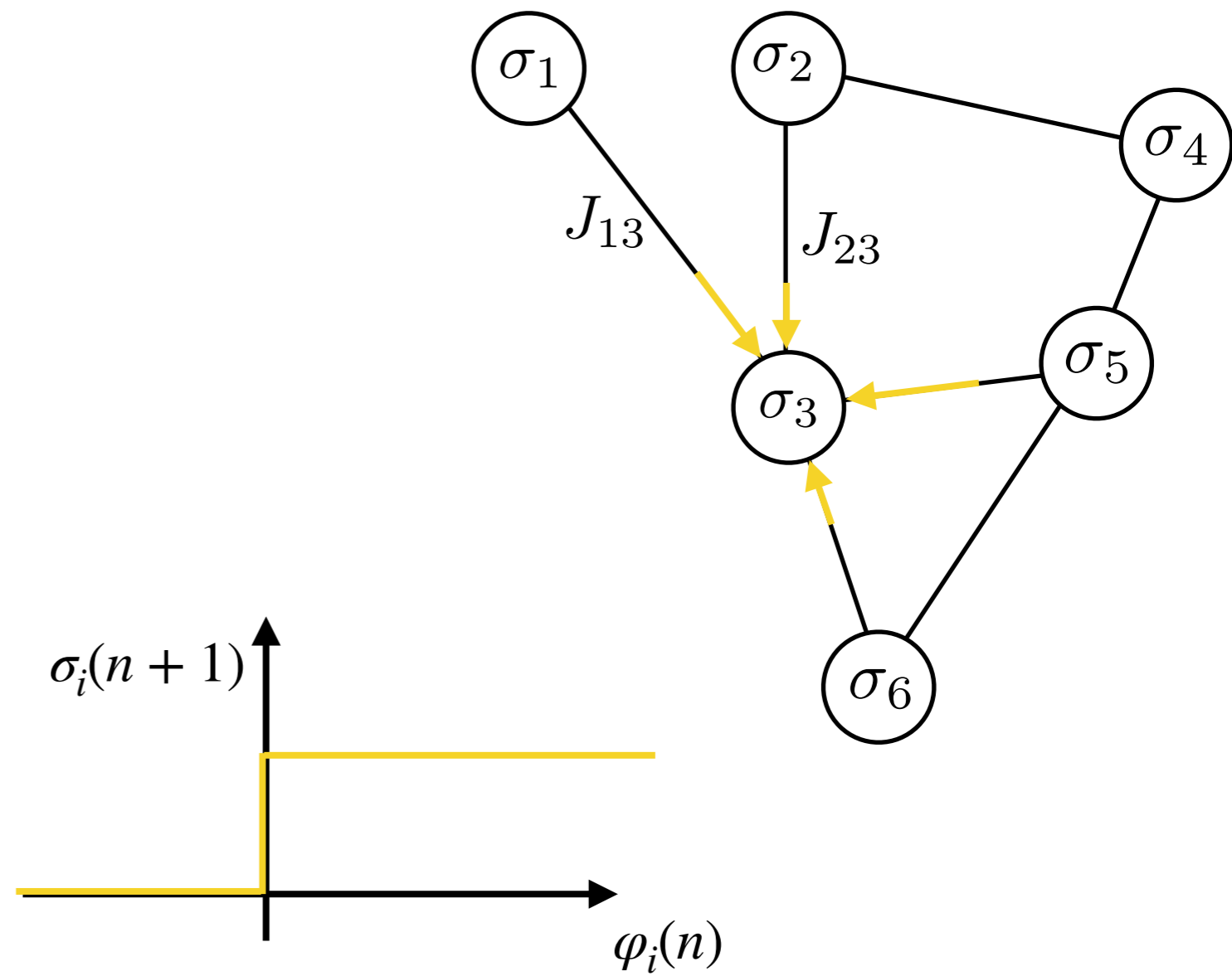


Neural dynamics

$$\sigma_i(n + 1) = \text{sgn}(\varphi_i(n))$$

Field acting on i
at time step n

$$\varphi_i(n) = \sum_{\substack{k=1 \\ k \neq i}}^N J_{ik} \sigma_k(n)$$



Neural dynamics

$$\sigma_i(n+1) = \text{sgn}(\varphi_i(n) + \zeta_i T)$$

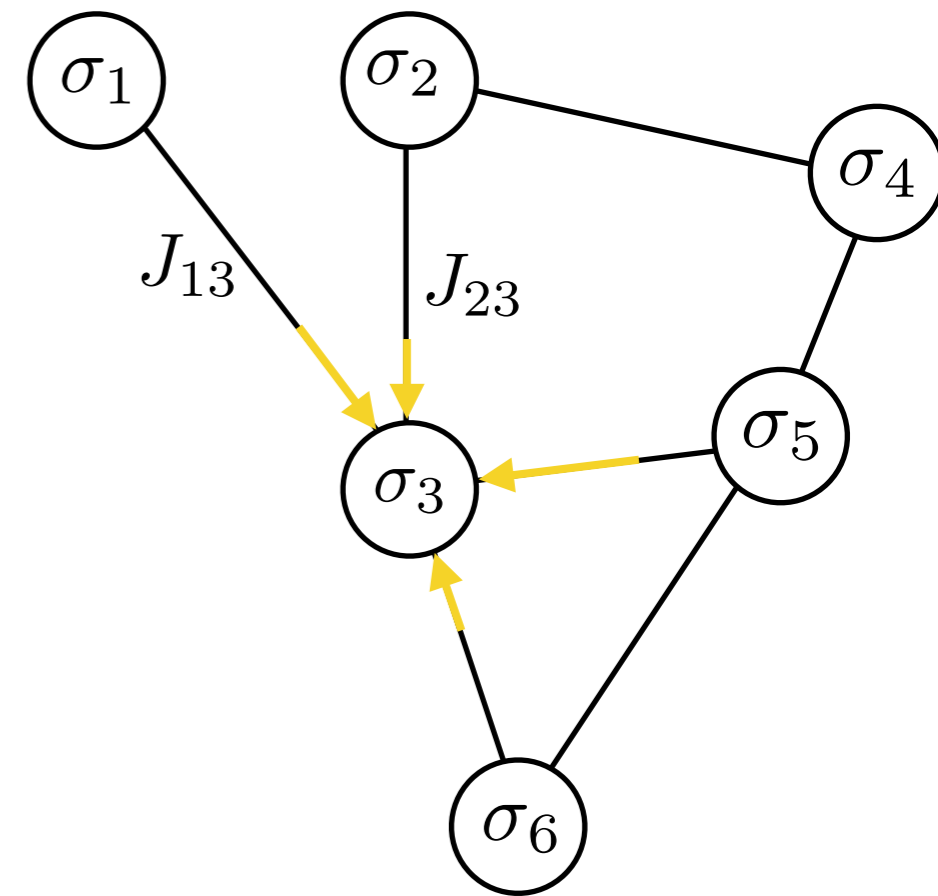
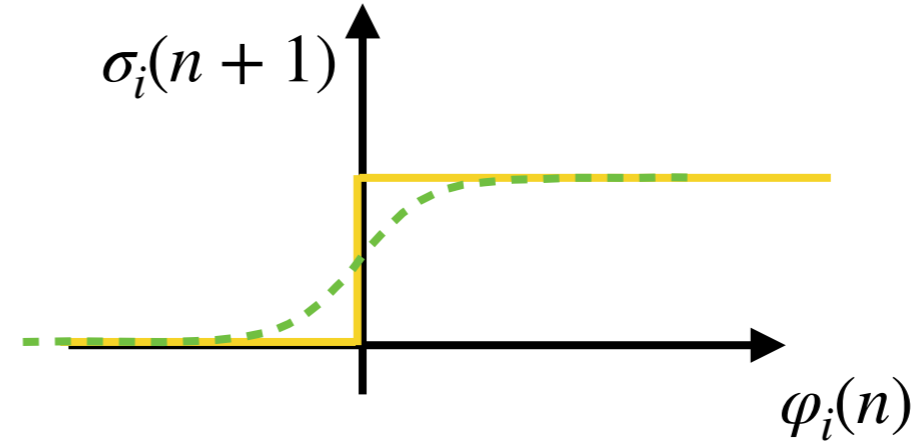
Field acting on i
at time step n

$$\varphi_i(n) = \sum_{\substack{k=1 \\ k \neq i}}^N J_{ik} \sigma_k(n)$$

Stochasticity

ζ_i i.i.d. r.v.'s

$T \geq 0$



Neural dynamics

$$\sigma_i(n+1) = \text{sgn}(\varphi_i(n) + \zeta_i T)$$

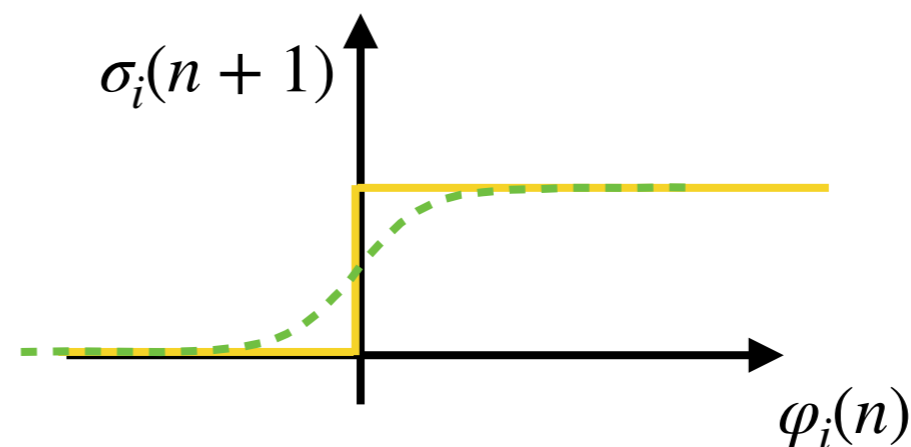
Field acting on i
at time step n

Stochasticity

ζ_i i.i.d. r.v.'s

$T \geq 0$

$$\varphi_i(n) = \sum_{\substack{k=1 \\ k \neq i}}^N J_{ik} \sigma_k(n)$$



If $T = 0$, as long as \mathbf{J} symmetric,

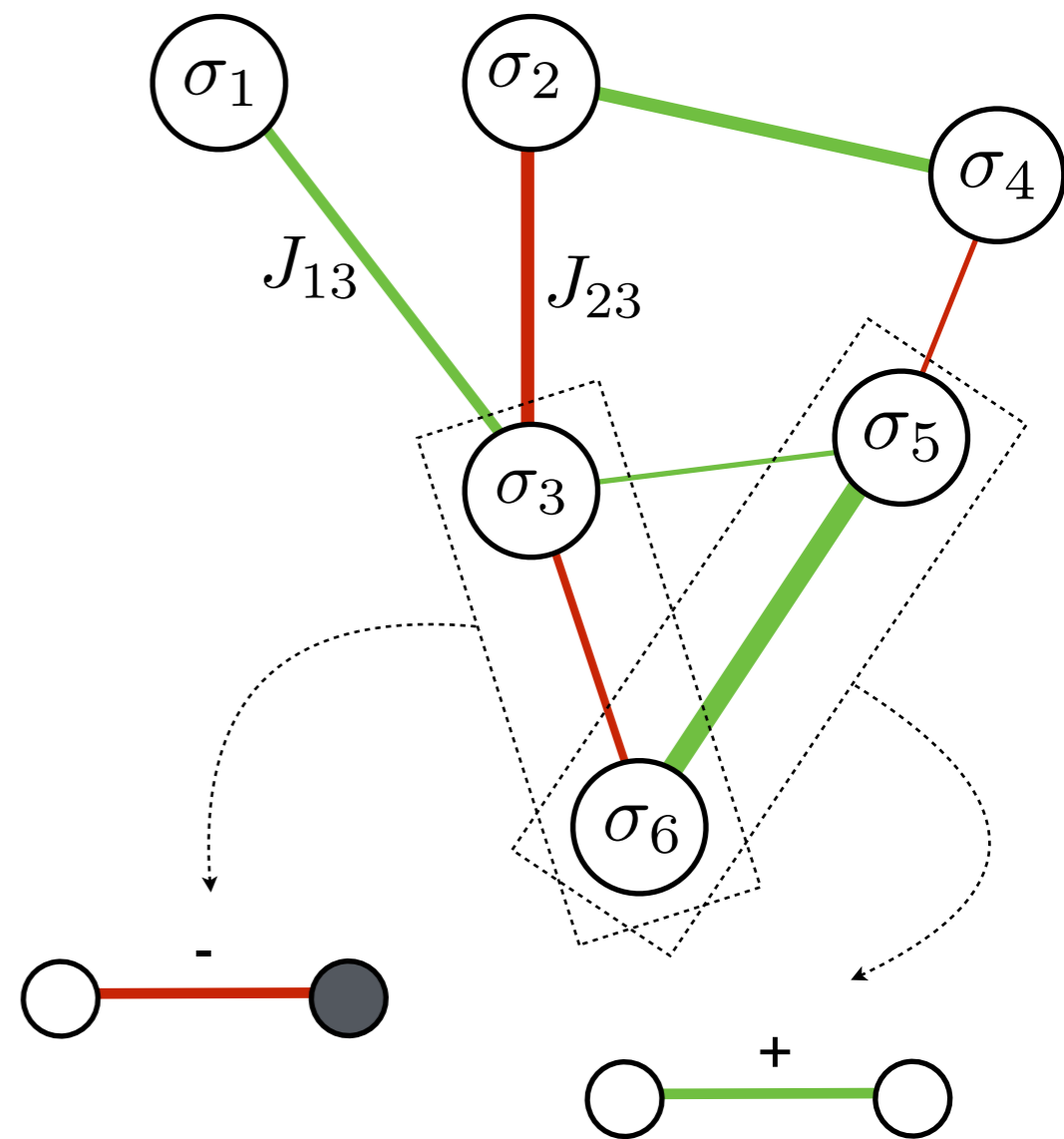
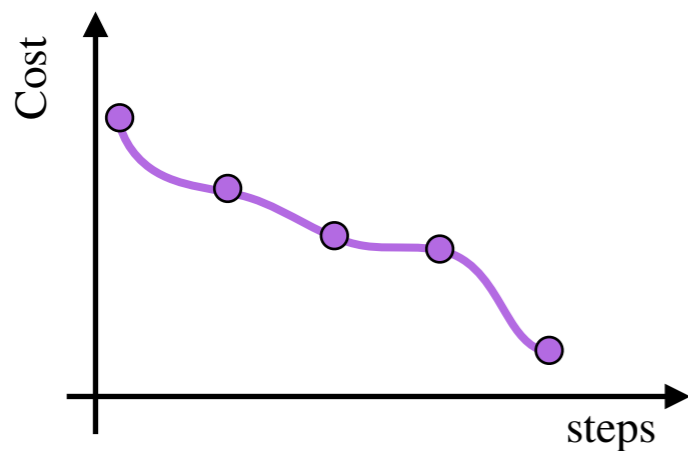
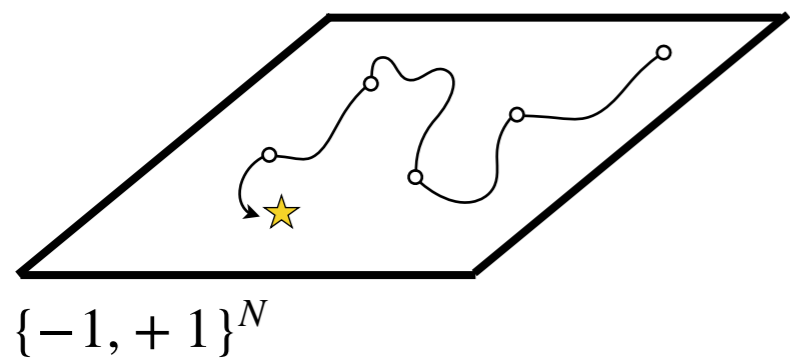
$$\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i,j} \sigma_i J_{ij} \sigma_j \text{ is a Lyapunov function}$$

If $T > 0$, as long as \mathbf{J} symmetric and $\text{Cum}(\zeta) = [1 + \tanh(\zeta)]/2$,

$$\mathcal{P}_{N,\beta,\mathbf{J}}(\boldsymbol{\sigma}) \propto e^{-\beta \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})}, \beta := T^{-1}$$

If $T = 0$, as long as \mathbf{J} symmetric,

$$\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i,j} \sigma_i J_{ij} \sigma_j \text{ is a Lyapunov function}$$



Fixed point dynamics

$$\boldsymbol{\sigma}^* = \operatorname{argmin} \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})$$

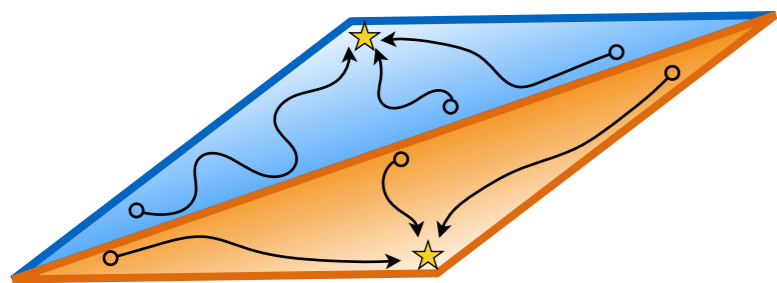
$$\boldsymbol{\sigma}_0 \rightarrow \boldsymbol{\sigma}^*(\boldsymbol{\sigma}_0, \mathbf{J})$$



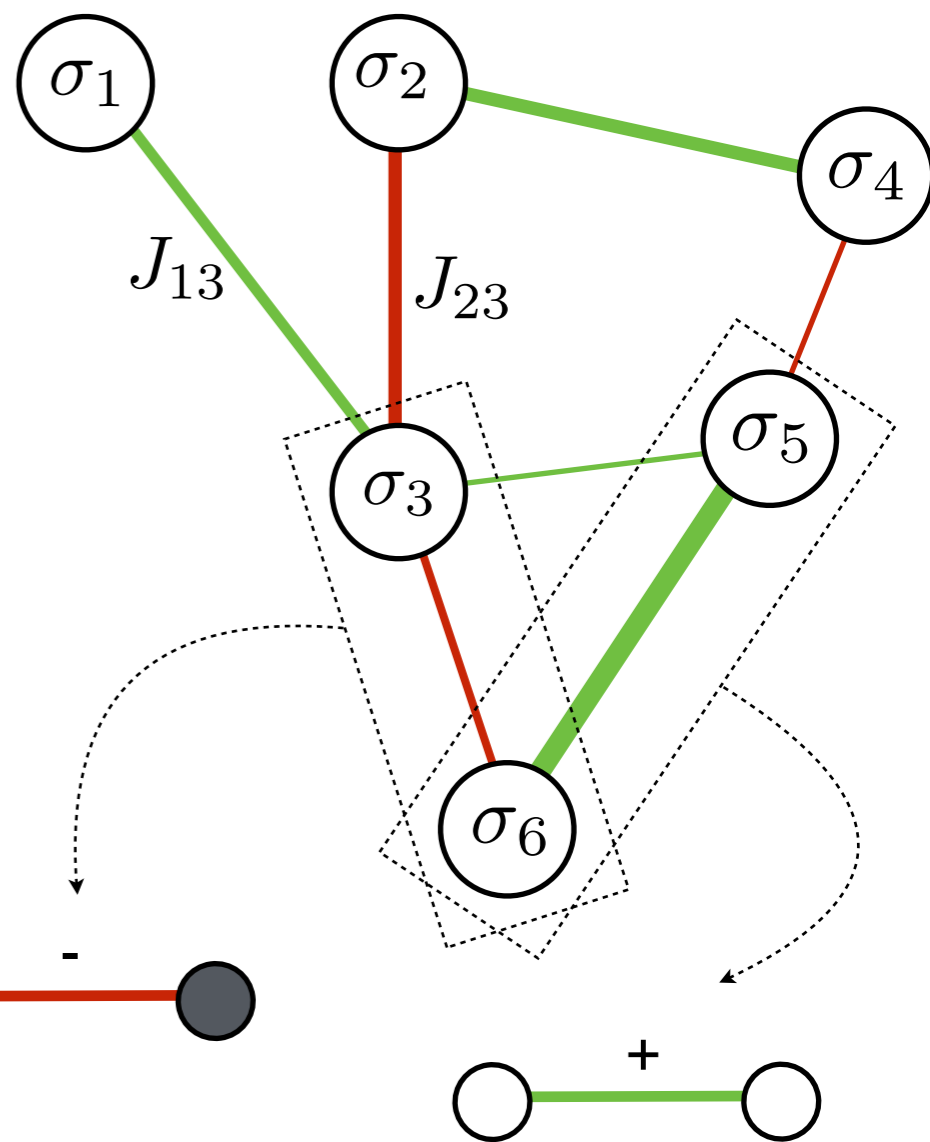
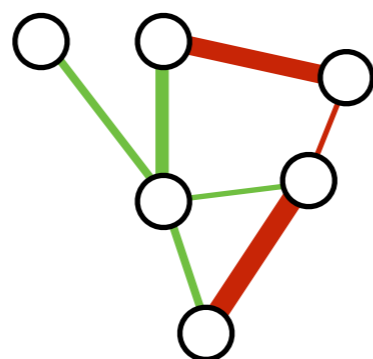
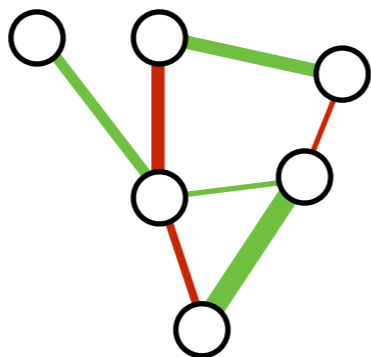
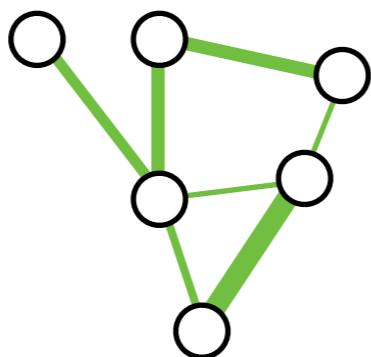
If $T = 0$, as long as \mathbf{J} symmetric,

$$\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i,j} \sigma_i J_{ij} \sigma_j \text{ is a Lyapunov function}$$

$$\boldsymbol{\sigma} = (-1, -1, \dots, -1)$$



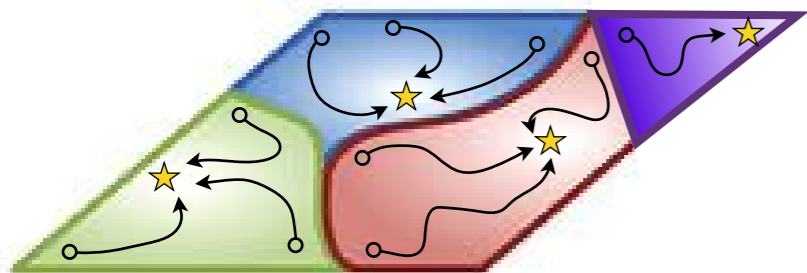
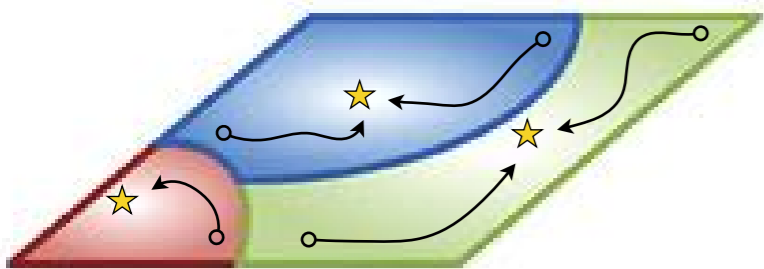
$$\boldsymbol{\sigma} = (+1, +1, \dots, +1)$$



Fixed point dynamics

$$\boldsymbol{\sigma}^* = \operatorname{argmin} \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})$$

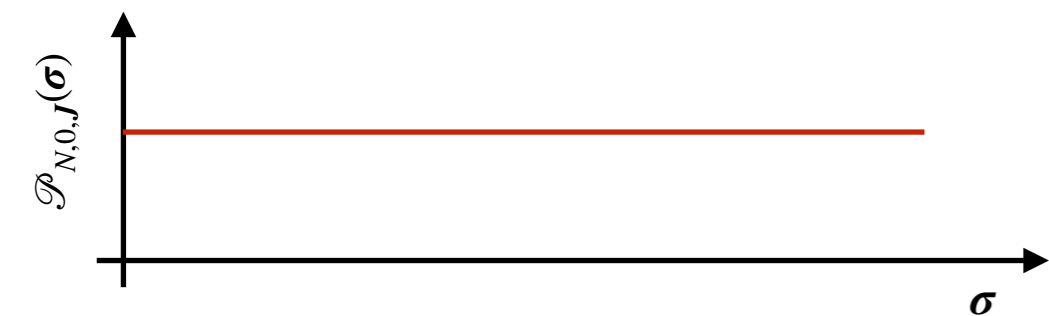
$$\boldsymbol{\sigma}_0 \rightarrow \boldsymbol{\sigma}^*(\boldsymbol{\sigma}_0, \mathbf{J})$$



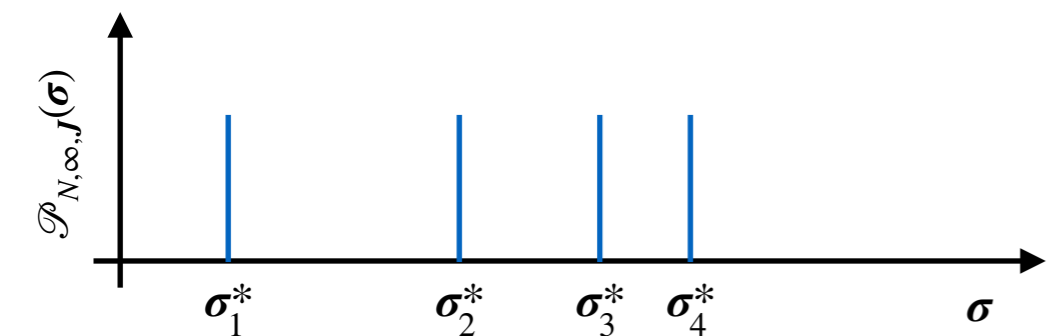
$$\sigma_i(n+1) = \text{sgn}(\varphi_i(n) + \zeta_i T)$$

If $T > 0$, as long as \mathbf{J} symmetric and $\text{Cum}(\zeta) = [1 + \tanh(\zeta)]/2$,

$$\mathcal{P}_{N,\beta,\mathbf{J}}(\boldsymbol{\sigma}) \propto e^{-\beta \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})}, \quad \beta := T^{-1}$$



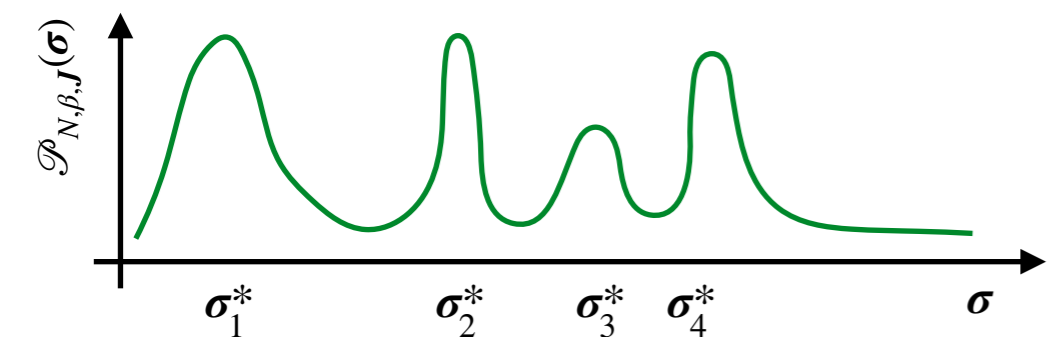
- Noisy limit $T \rightarrow \infty$
Fully random dynamics $\sigma_i(n+1) = \pm 1$



- Noiseless limit $T \rightarrow 0$
Deterministic dynamics $\sigma_i(n+1) = \text{sgn}(\varphi_i(n))$
 $\mathcal{P}_{N,\beta,\mathbf{J}}(\boldsymbol{\sigma})$ delta-peaked at $\boldsymbol{\sigma}^*$
Fixed points dynamics

$$\boldsymbol{\sigma}^* = \text{argmin } \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})$$

$$\boldsymbol{\sigma}_0 \rightarrow \boldsymbol{\sigma}^*(\boldsymbol{\sigma}_0, \mathbf{J})$$



- Intermediate regime $T \in \mathbb{R}^+$
 $\mathcal{P}_{N,T,\mathbf{J}}(\boldsymbol{\sigma})$ peaked at $\boldsymbol{\sigma}^*$
Starting at $\boldsymbol{\sigma}_0$, likely to end up in neighbourhood of $\boldsymbol{\sigma}^*(\boldsymbol{\sigma}_0, \mathbf{J})$

The statistical mechanics setting

Mean-field spin-like model: $\sigma \in \{-1, +1\}^N$, \mathbf{J} symm

Hamiltonian $\mathcal{H}_{N,\mathbf{J}}(\sigma) = -\frac{1}{2} \sum_{i,j=1}^N \sigma_i J_{ij} \sigma_j$

Boltzmann-Gibbs $\mathcal{P}_{N,\beta,\mathbf{J}}(\sigma) = \frac{e^{-\beta \mathcal{H}_{N,\mathbf{J}}(\sigma)}}{\mathcal{Z}_{N,\beta,\mathbf{J}}}$, with $\mathcal{Z}_{N,\beta,\mathbf{J}} = \sum_{\{\sigma\}} \exp[-\beta \mathcal{H}_{N,\mathbf{J}}(\sigma)]$, $\beta := T^{-1}$

The statistical mechanics setting

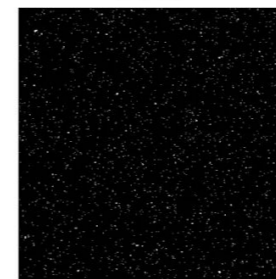
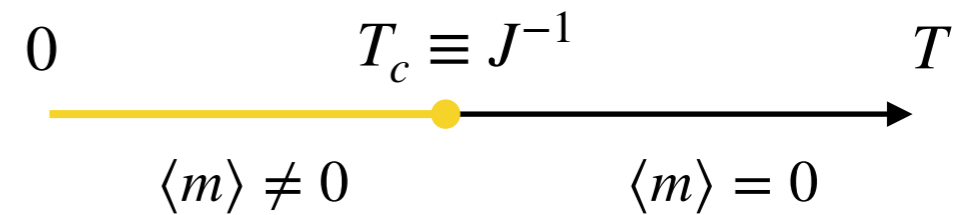
Mean-field spin-like model: $\sigma \in \{-1, +1\}^N$, \mathbf{J} symm

Hamiltonian $\mathcal{H}_{N,\mathbf{J}}(\sigma) = -\frac{1}{2} \sum_{i,j=1}^N \sigma_i J_{ij} \sigma_j$

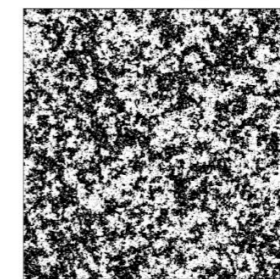
Boltzmann-Gibbs $\mathcal{P}_{N,\beta,\mathbf{J}}(\sigma) = \frac{e^{-\beta \mathcal{H}_{N,\mathbf{J}}(\sigma)}}{\mathcal{Z}_{N,\beta,\mathbf{J}}}$, with $\mathcal{Z}_{N,\beta,\mathbf{J}} = \sum_{\{\sigma\}} \exp[-\beta \mathcal{H}_{N,\mathbf{J}}(\sigma)]$, $\beta := T^{-1}$

Curie-Weiss model: $J_{ij} = J > 0, \forall_{i,j}$
(ferromagnet)

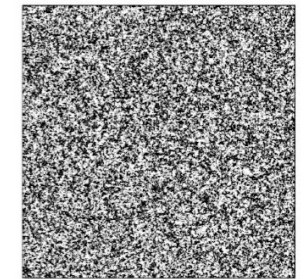
$$m_N(\sigma) := \frac{1}{N} \sum_{i=1}^N \sigma_i$$



$T < T_c$



$T \sim T_c$



$T > T_c$

The statistical mechanics setting

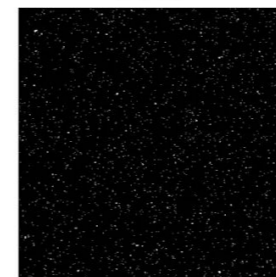
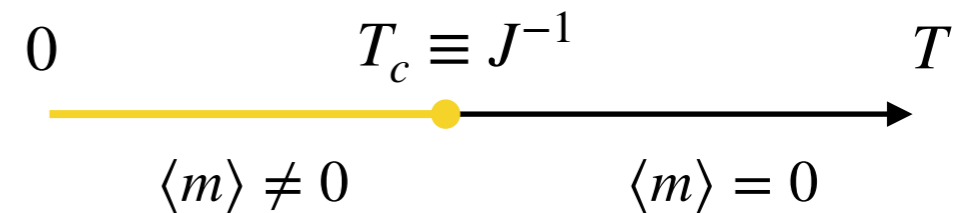
Mean-field spin-like model: $\sigma \in \{-1, +1\}^N$, \mathbf{J} symm

Hamiltonian $\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i,j=1}^N \sigma_i J_{ij} \sigma_j$

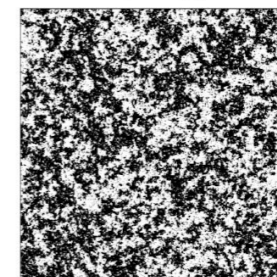
Boltzmann-Gibbs $\mathcal{P}_{N,\beta,\mathbf{J}}(\boldsymbol{\sigma}) = \frac{e^{-\beta \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})}}{\mathcal{Z}_{N,\beta,\mathbf{J}}}$, with $\mathcal{Z}_{N,\beta,\mathbf{J}} = \sum_{\{\boldsymbol{\sigma}\}} \exp[-\beta \mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma})]$, $\beta := T^{-1}$

Curie-Weiss model: $J_{ij} = J > 0, \forall_{i,j}$
(ferromagnet)

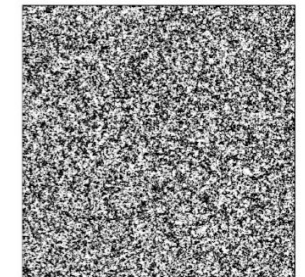
$$m_N(\boldsymbol{\sigma}) := \frac{1}{N} \sum_{i=1}^N \sigma_i$$



$T < T_c$



$T \sim T_c$



$T > T_c$

Find explicit expression for *free-energy* $\mathcal{F} := -T \log \mathcal{Z}$

→ n -th moment of macroscopic observables related to free-energy n -th order derivatives w.r.t. conjugate parameters

Neural networks for retrieval

Mimics *retrieval* capabilities

Pattern recognition

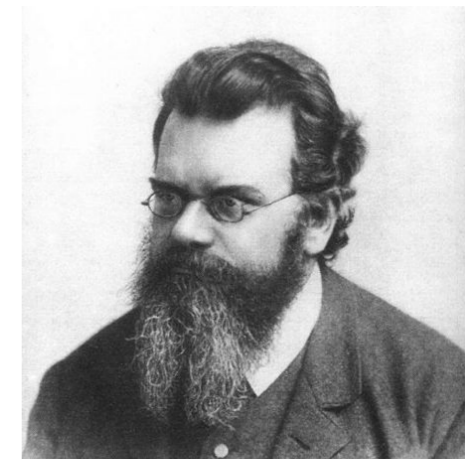
Pattern reconstruction

Denoising

Content addressable memory



CAPTCHA



Neural networks for retrieval

Mimics *retrieval* capabilities

Pattern recognition

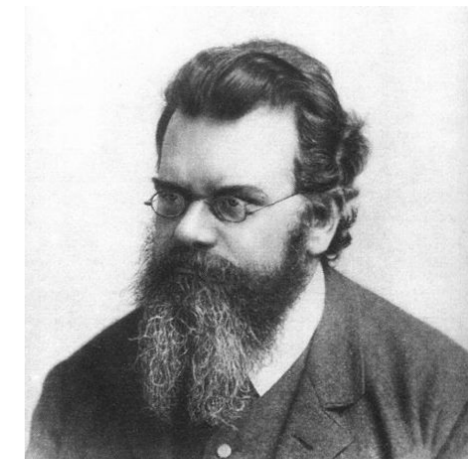
Pattern reconstruction

Denoising

Content addressable memory

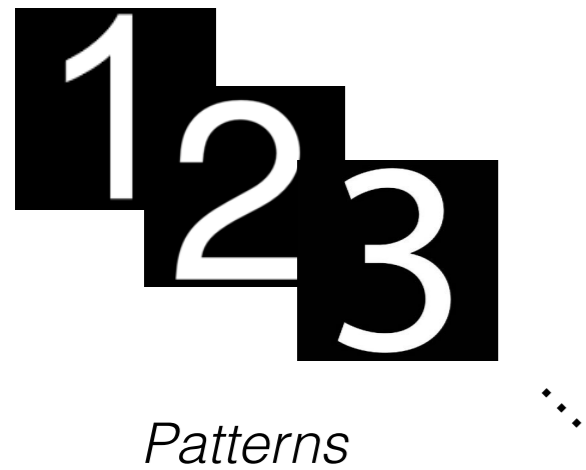


CAPTCHA



Network asked to reconstruct previously learnt vectors, now inputted noisy or incomplete

Set of definite K patterns of length N binary vectors (archetypes)



$$\xi^1, \dots, \xi^K$$

$$\xi^\mu = (+1, -1, \dots, +1) \in \{-1, +1\}^N$$

input : σ_0

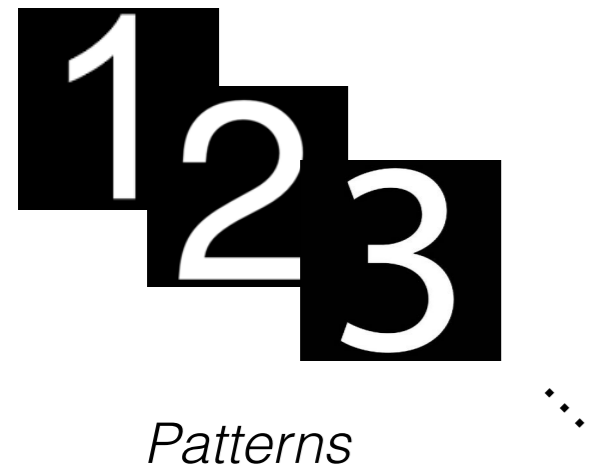


output : $\sigma^*(\sigma_0)$



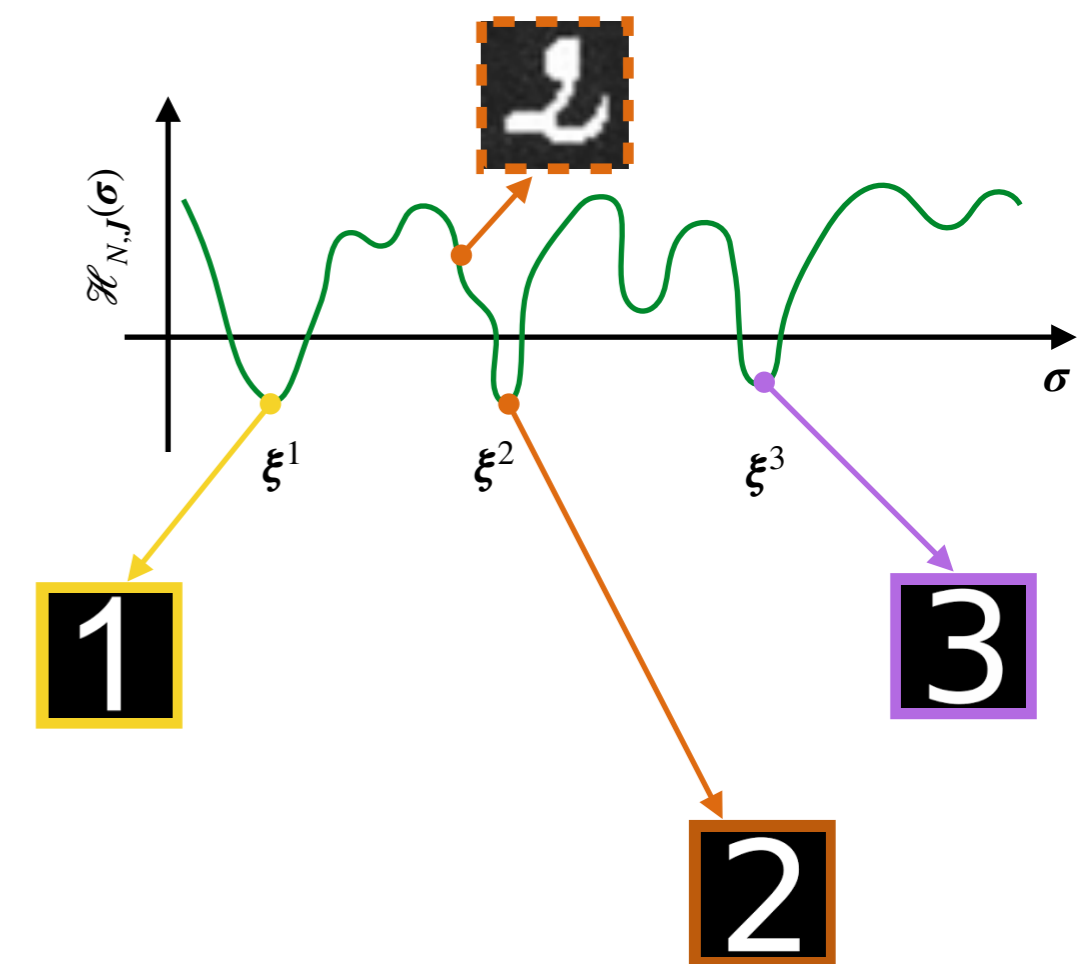
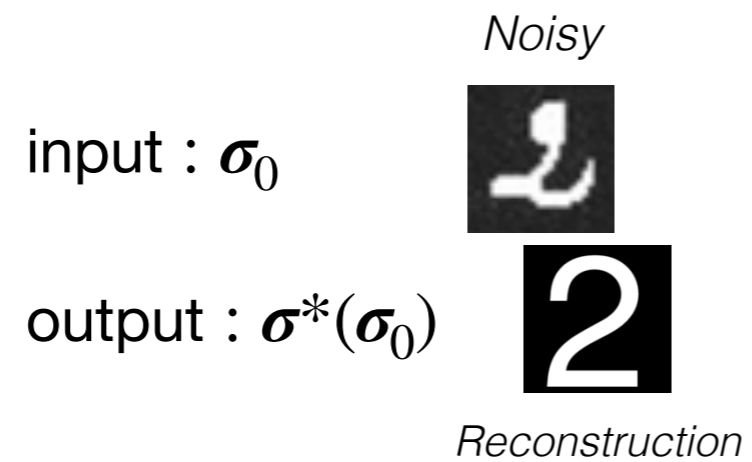
Network asked to reconstruct previously learnt vectors, now inputted noisy or incomplete

Set of definite K patterns of length N binary vectors (archetypes)



$$\xi^1, \dots, \xi^K$$

$$\xi^\mu = (+1, -1, \dots, +1) \in \{-1, +1\}^N$$

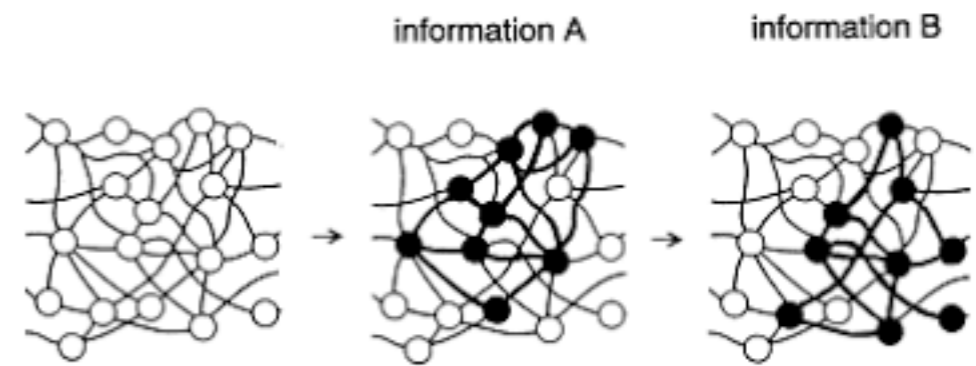


$$\sigma(n=0) = \sigma_0 \rightarrow \sigma^*(\sigma_0) = \lim_{n \rightarrow \infty} \sigma(n) \Big|_{\sigma_0}$$

Choose $J = J(\{\xi_i^\mu\}_{i=1, \dots, N}^{\mu=1, \dots, K})$ s.t. $\xi^\mu = \operatorname{argmin}_{\sigma \in \{-1, +1\}^N} \mathcal{H}_{N,J}(\sigma)$ for $\mu = 1, \dots, K$

The Hopfield model (Hopfield model '82, Pastur&Figotin '78)

$$\text{Hebb's rule } J_{ij} = \frac{1}{N} \sum_{\mu=1}^K \xi_i^{\mu} \xi_j^{\mu}$$

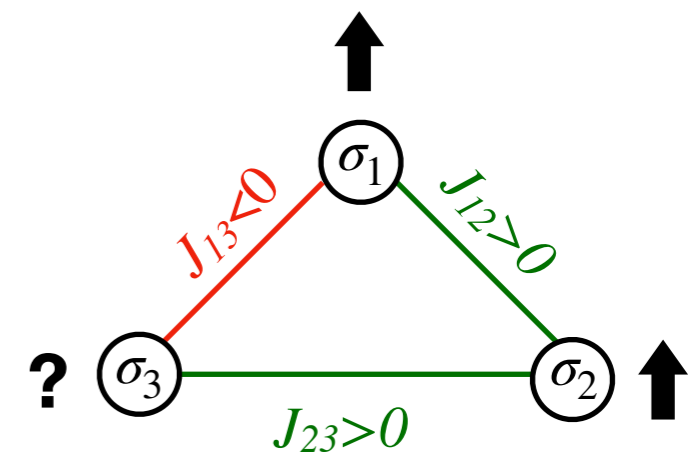
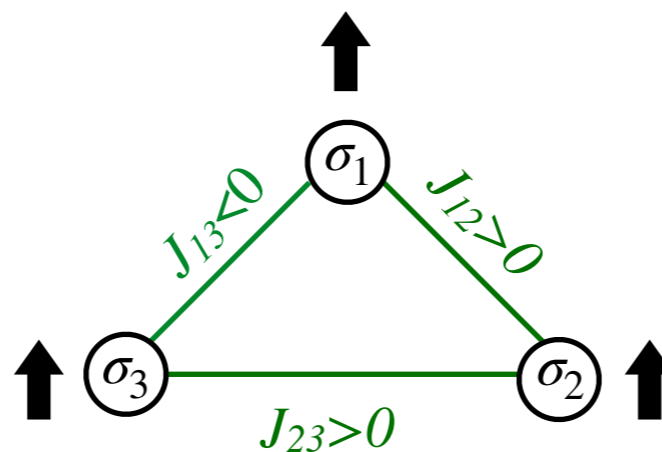


“Neurons out of sync fail to link”

- Biological inspiration
- Low complexity
- Sub-optimal
- Excitatory vs Inhibitory

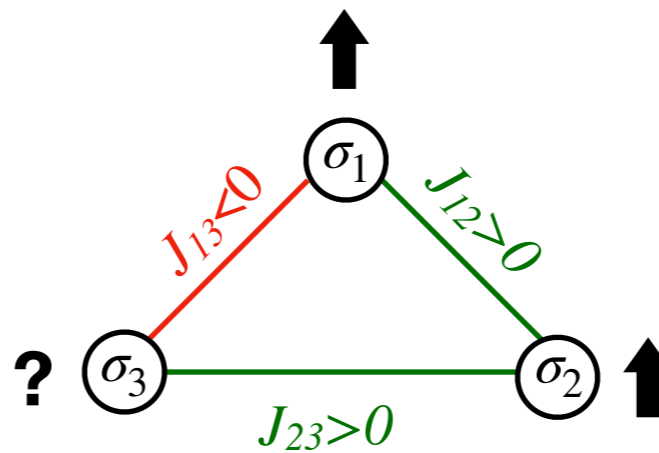
$$\mathcal{H}_{N,J}(\boldsymbol{\sigma}) = - \sum_{i>j} \sigma_i J_{ij} \sigma_j$$

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^K \xi_i^{\mu} \xi_j^{\mu}$$



$$\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = - \sum_{i>j} \sigma_i J_{ij} \sigma_j$$

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu$$



Frustration makes the energy landscape rough with #local-minima exponentially increasing with N

NP-complete combinatorial optimization problem

Load $\alpha := \lim_{N \rightarrow \infty} \frac{K}{N}, \alpha > 0$

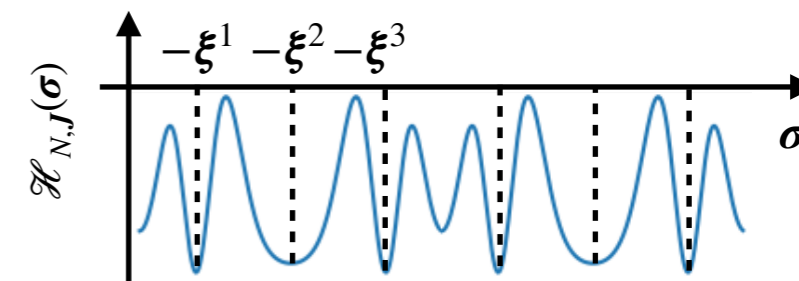
trivial

$$K = 0$$



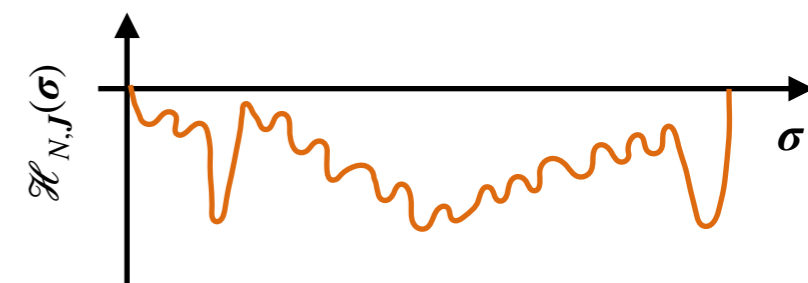
simple

$$K \sim \mathcal{O}(1)$$



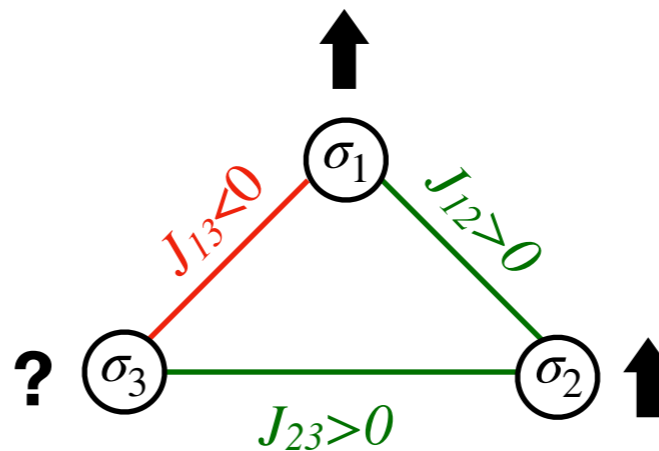
complex

$$K \sim \mathcal{O}(N)$$



$$\mathcal{H}_{N,\mathbf{J}}(\boldsymbol{\sigma}) = - \sum_{i>j} \sigma_i J_{ij} \sigma_j$$

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu$$



Frustration makes the energy landscape rough with #local-minima exponentially increasing with N

NP-complete combinatorial optimization problem

Load $\alpha := \lim_{N \rightarrow \infty} \frac{K}{N}, \alpha > 0$

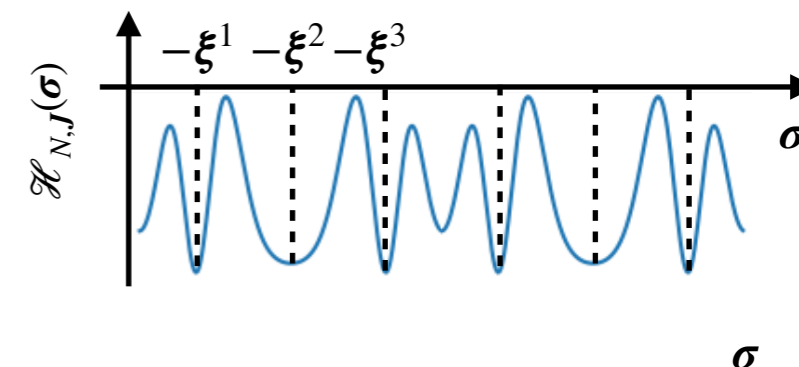
trivial

$$K = 0$$



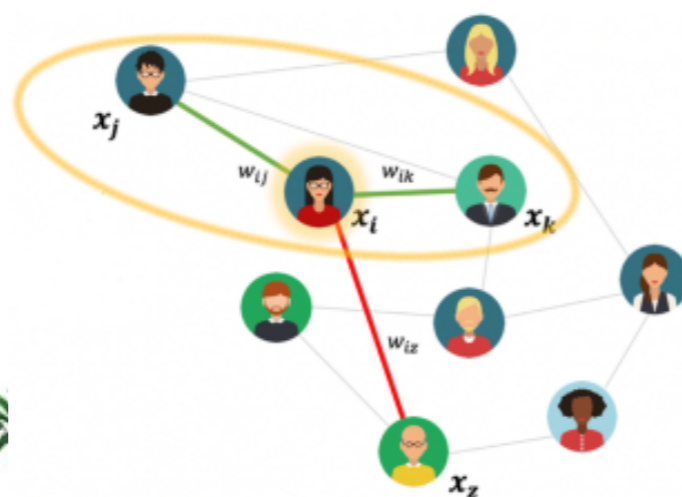
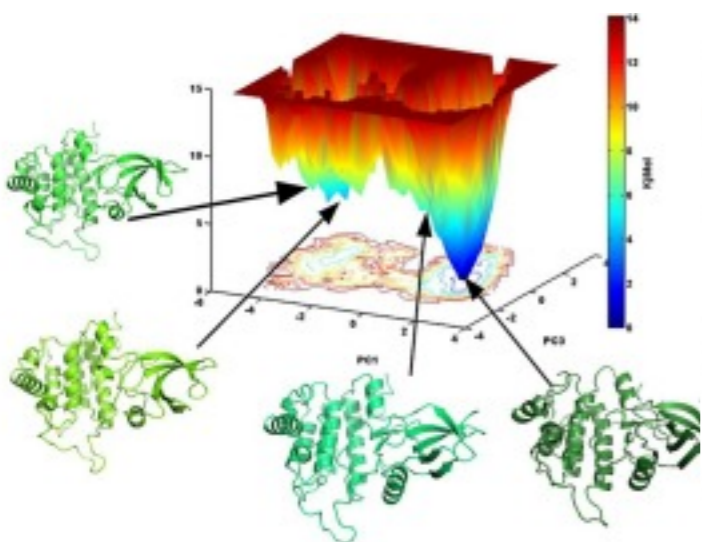
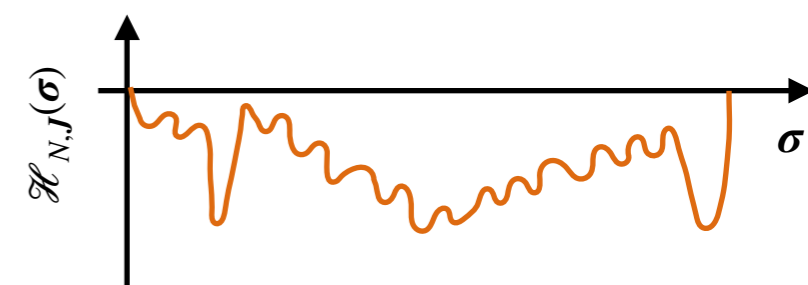
simple

$$K \sim \mathcal{O}(1)$$

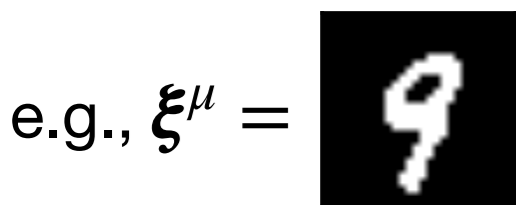


complex

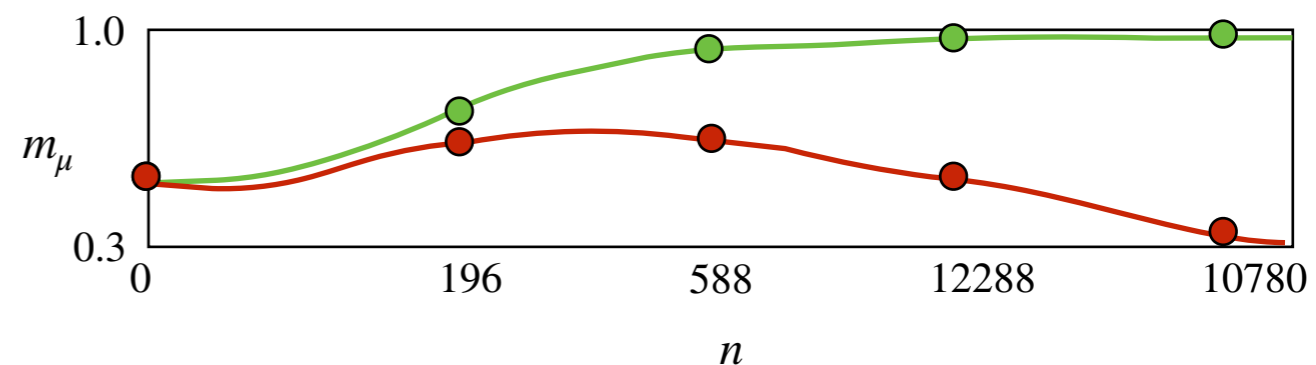
$$K \sim \mathcal{O}(N)$$



The model *works* (as a pattern reconstructor) if starting “near” a pattern, say $\sigma \approx \xi^\mu$, the dynamics eventually converges to that pattern. Here “near” means that the Hamming distance between the initial point and the pattern is smaller than ϵN , for some $\epsilon > 0$.



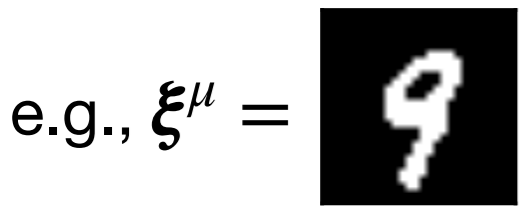
$n = 0$ $n = 196$ $n = 588$ $n = 2940$ $n = 10780$



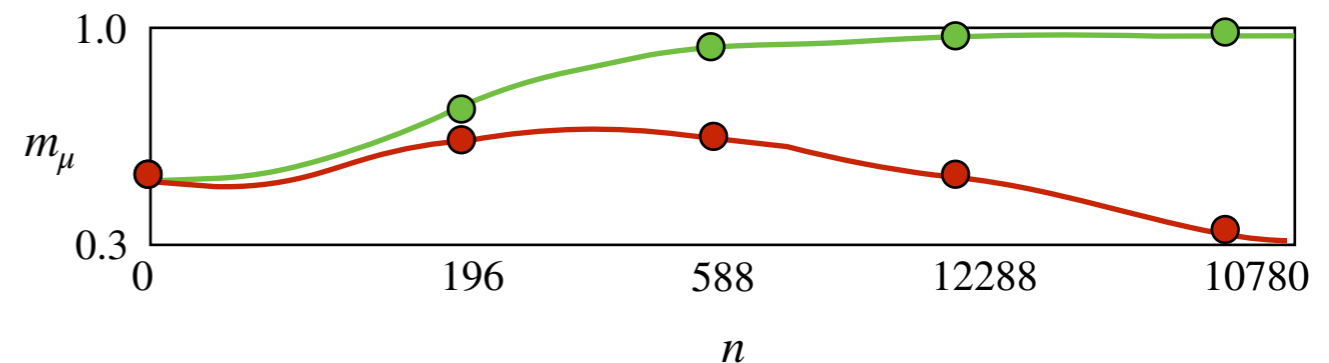
Mattis magnetization $m_\mu := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i$

assess retrieval of the μ -th pattern

The model *works* (as a pattern reconstructor) if starting “near” a pattern, say $\sigma \approx \xi^\mu$, the dynamics eventually converges to that pattern. Here “near” means that the Hamming distance between the initial point and the pattern is smaller than ϵN , for some $\epsilon > 0$.

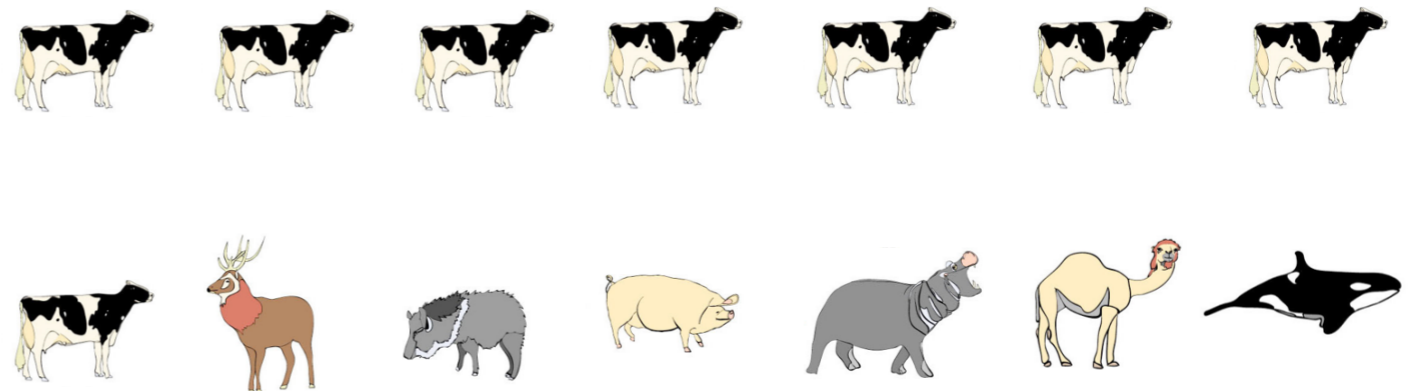


$n = 0$ $n = 196$ $n = 588$ $n = 2940$ $n = 10780$



Mattis magnetization $m_\mu := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i$
 assess retrieval of the μ -th pattern

Overlap $q_{ab} := \frac{1}{N} \sum_{i=1}^N \sigma_i^{(a)} \sigma_i^{(b)}$



intrinsic disorder and frustration yield glassy behaviours, captured by the correlation between configurations of two *replicas*

AGS theory

$$\mathcal{H}_{N,\xi}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{\substack{i,j=1 \\ i \neq j}}^N \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j$$

M Mattis magnetization $m_\mu := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i$,

assess retrieval of the μ -th pattern

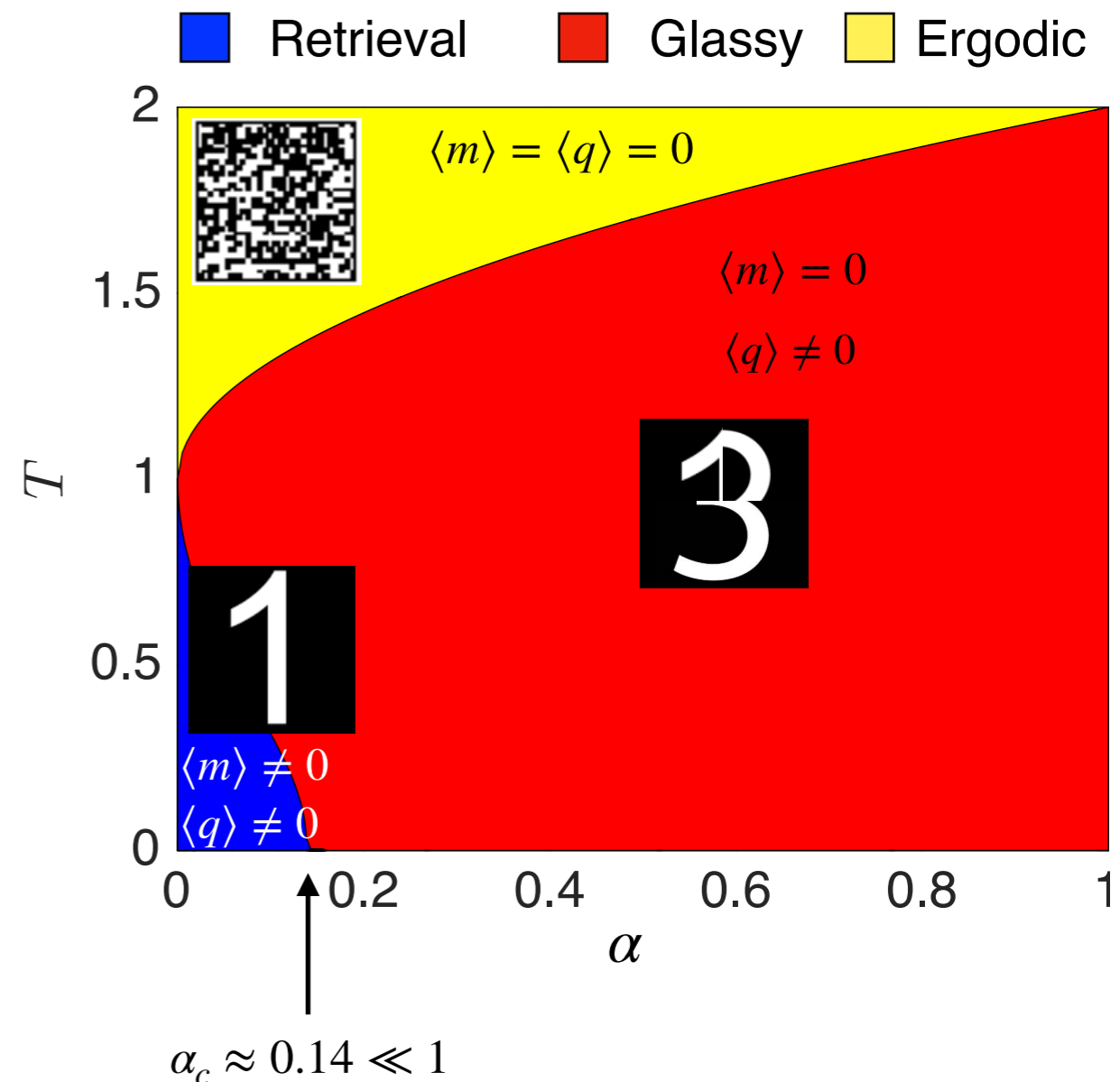
Overlap $q_{ab} := \frac{1}{N} \sum_{i=1}^N \sigma_i^{(a)} \sigma_i^{(b)}$

describes correlation between configurations of two *replicas*, namely two different systems endowed with the same couplings \mathbf{J}

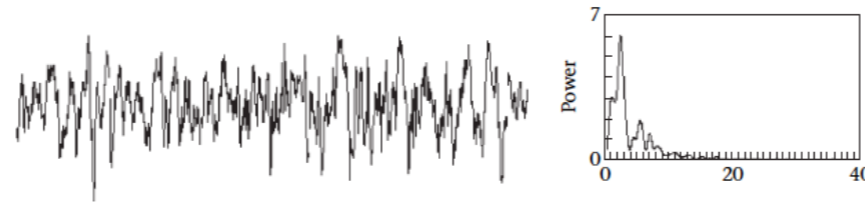
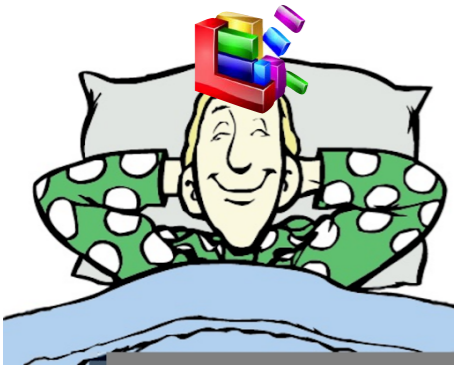
Load $\alpha := \lim_{N \rightarrow \infty} \frac{K}{N}$

Noise $T := \beta^{-1}$

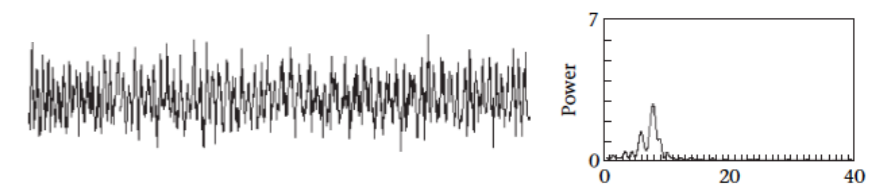
[Amit, Gutfreund, Sompolinsky - Phys. Rev. A '87]



Revise Hebb's rule by iterative protocols



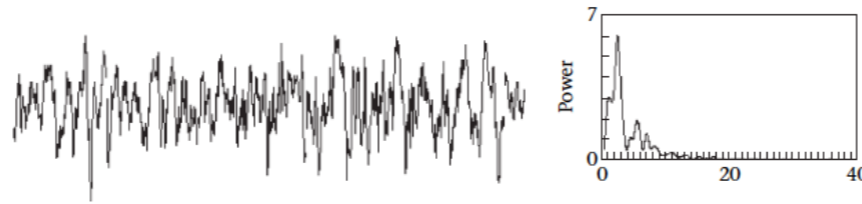
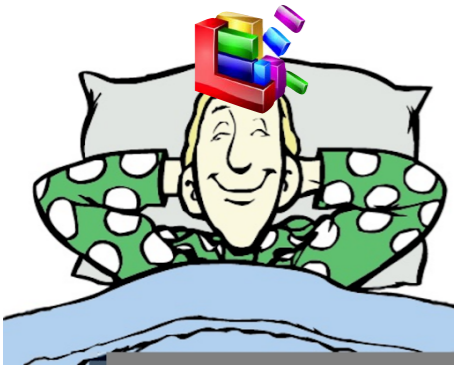
consolidation (SW)



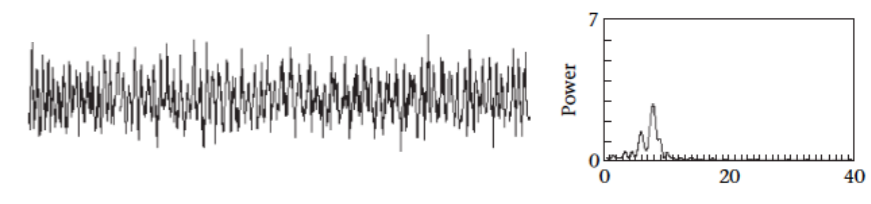
remotion (REM)

[Crick, Mitchinson - Nature '83; Stickgold - Nature '05; Diekelmann, Born - Nature Rev. Neurosc. '10]

Revise Hebb's rule by iterative protocols



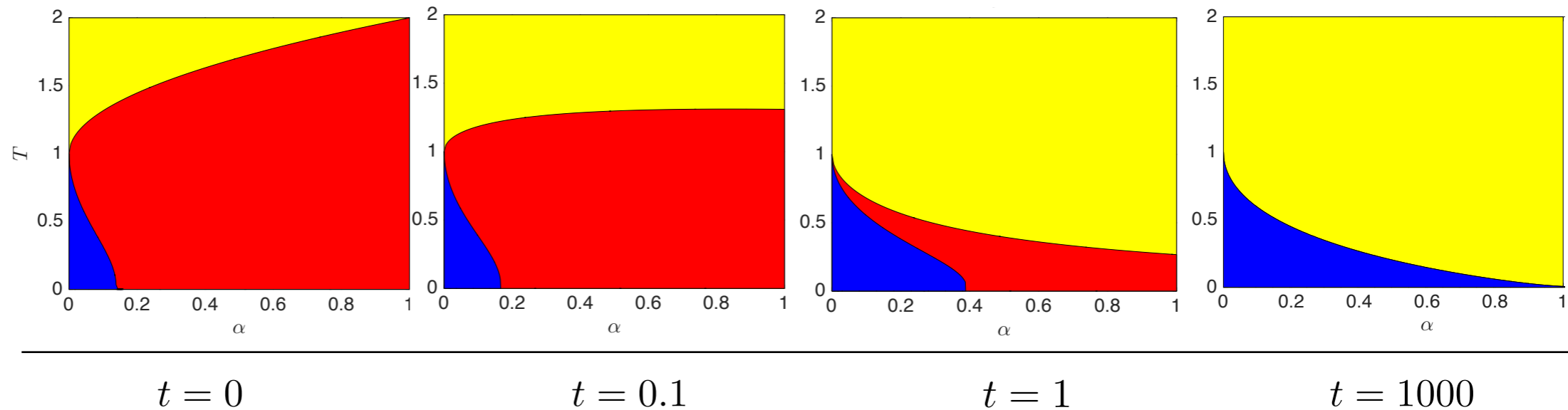
consolidation (SW)



remotion (REM)

[Crick, Mitchinson - Nature '83; Stickgold - Nature '05; Diekelmann, Born - Nature Rev. Neurosc. '10]

Retrieval ■
Spin-glass ■
Ergodic ■



Allocate more information with the same resources

$$J_{ij}^{(n)} \leftarrow J_{ij}^{(n-1)} + \frac{1}{1+n} [J_{ij}^{(n-1)} - (J_{ij}^{(n-1)})^2]$$

$$J_{ij}(t) = \frac{1}{N} \sum_{\mu, \nu} \xi_i^\mu \left(\frac{1+t}{1+t\mathbf{C}} \right)_{\mu\nu} \xi_j^\nu$$

$$\underbrace{(1+t)}_{\text{consolidation (SW)}} \cdot \underbrace{(\mathbf{I} + t\mathbf{C})^{-1}}_{\text{remotion (REM)}}$$

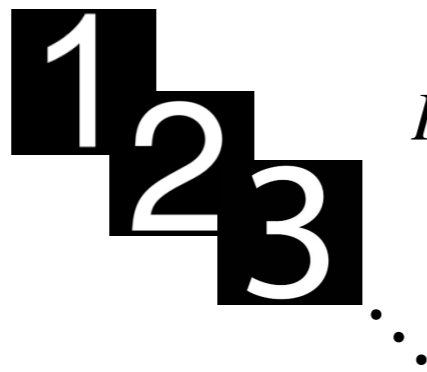
t "sleeping time"

where $C_{\mu\nu} := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \xi_i^\nu$ pattern correlation matrix

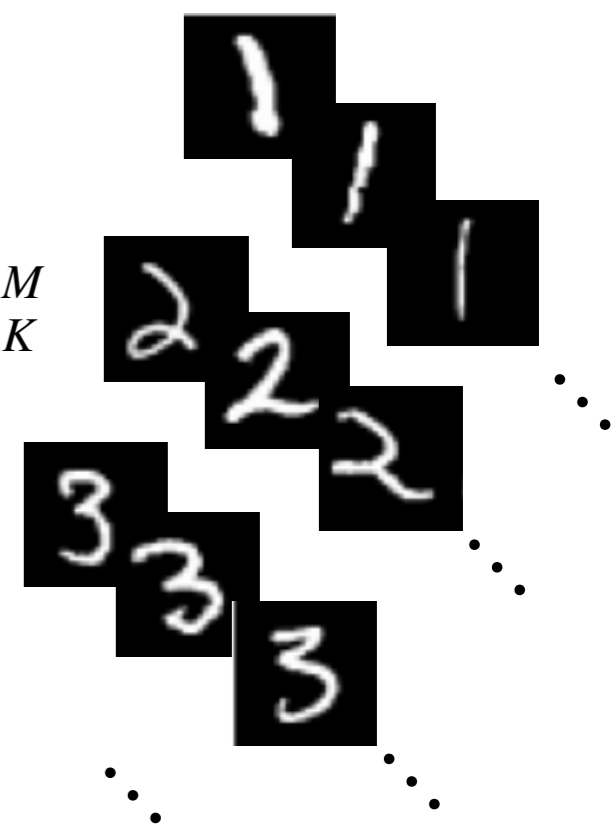
Patterns' revision

Replace *archetypes* by a sample of *examples*

K "archetypes" $\{\xi^\mu\}_{\mu=1,\dots,K}$



$K \times M$ "blurred" examples $\{\eta^{\mu,a}\}_{\mu=1,\dots,K}^{a=1,\dots,M}$



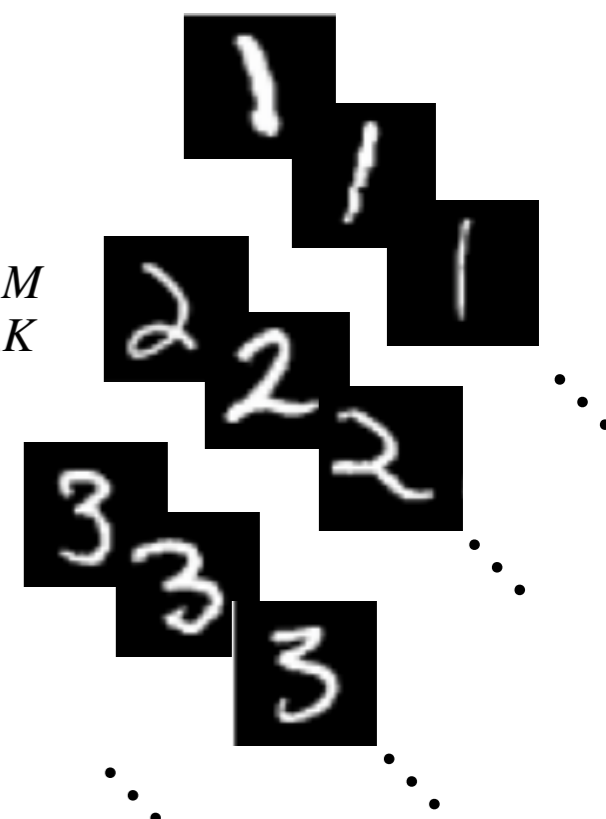
Patterns' revision

Replace *archetypes* by a sample of *examples*

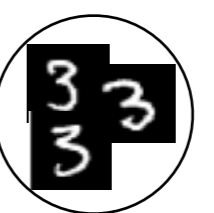
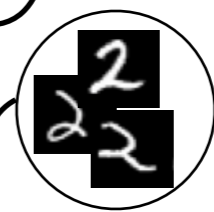
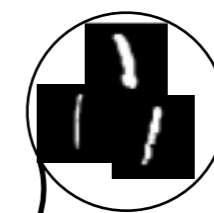
K "archetypes" $\{\xi^\mu\}_{\mu=1,\dots,K}$



$K \times M$ "blurred" examples $\{\eta^{\mu,a}\}_{\mu=1,\dots,K}^{a=1,\dots,M}$

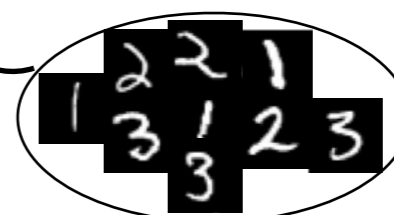


$$J_{ij}^{(sup)} \propto \sum_{\mu=1}^K \left(\sum_{a=1}^M \eta_i^{\mu,a} \right) \left(\sum_{b=1}^M \eta_i^{\mu,b} \right)$$



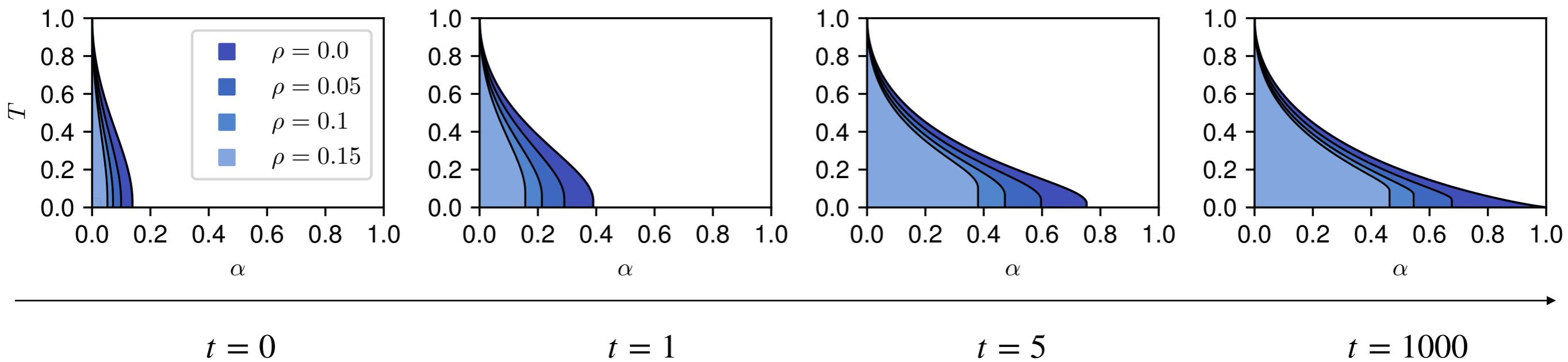
$$J_{ij}(\xi) \propto \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu$$

$$J_{ij}^{(unsup)} \propto \sum_{\mu=1}^K \sum_{a=1}^M \eta_i^{\mu,a} \eta_j^{\mu,a}$$



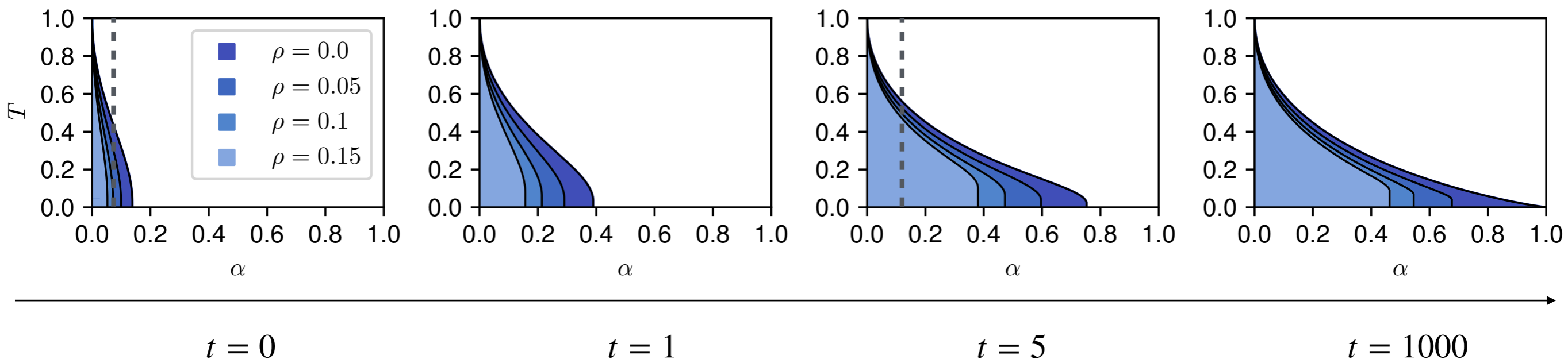
T, α, t

$$\rho \propto \frac{1}{\text{“quality”}} \frac{1}{\text{“quantity”}}$$

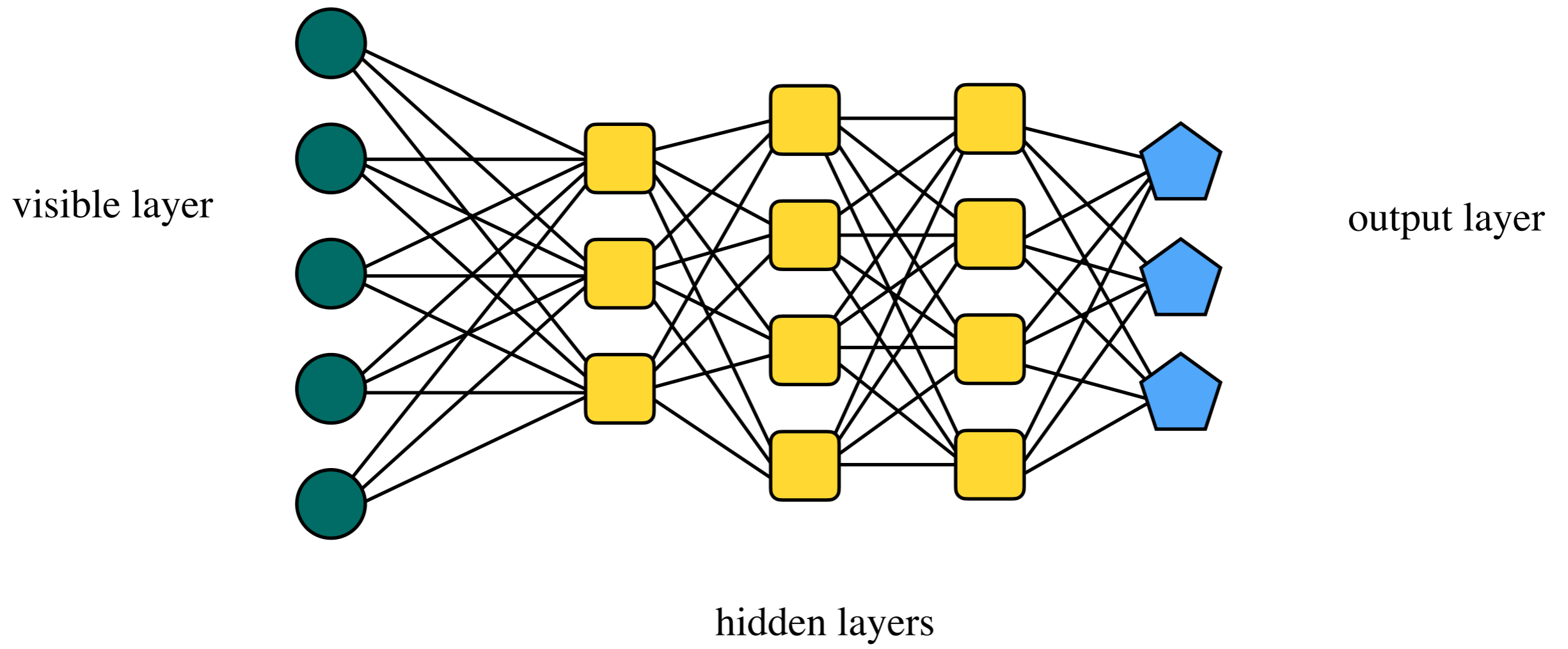


T, α, t

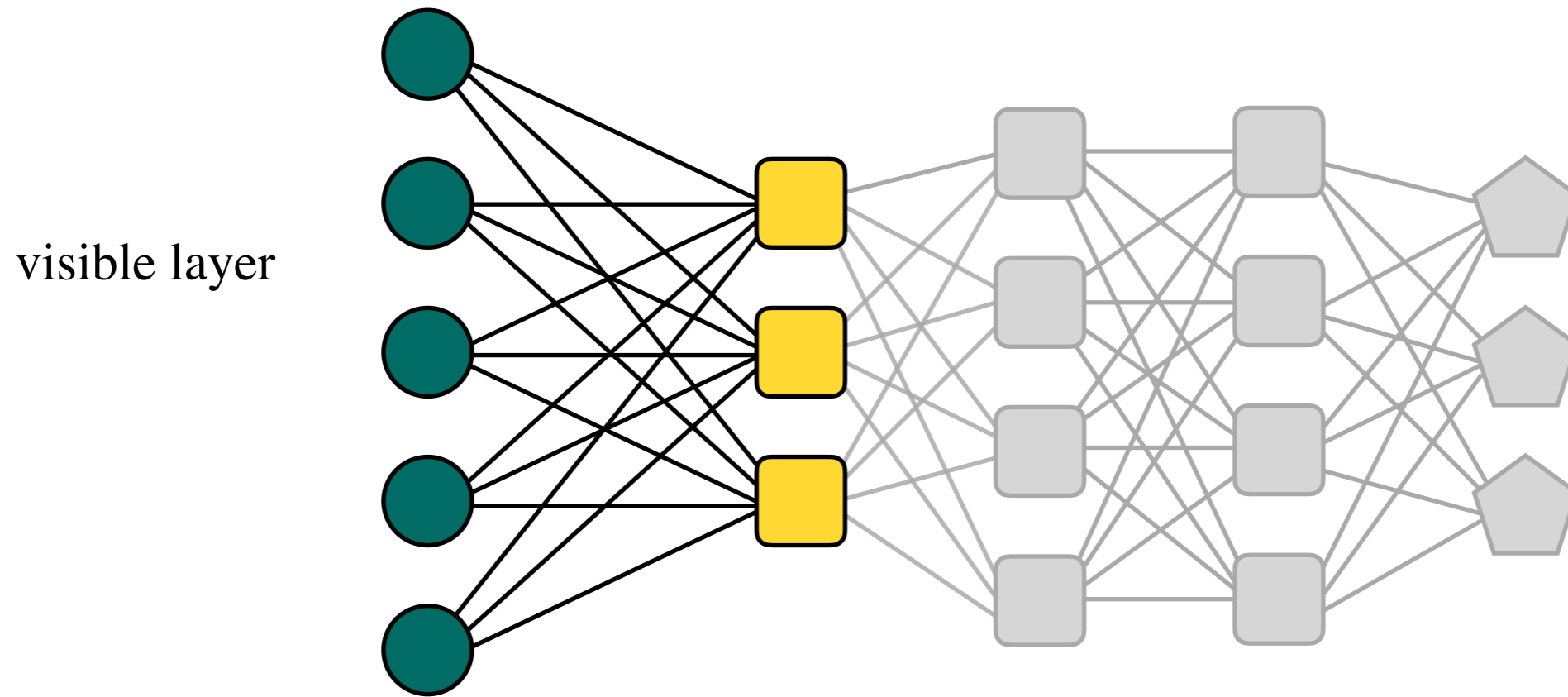
$$\rho \propto \frac{1}{\text{“quality”}} \frac{1}{\text{“quantity”}}$$



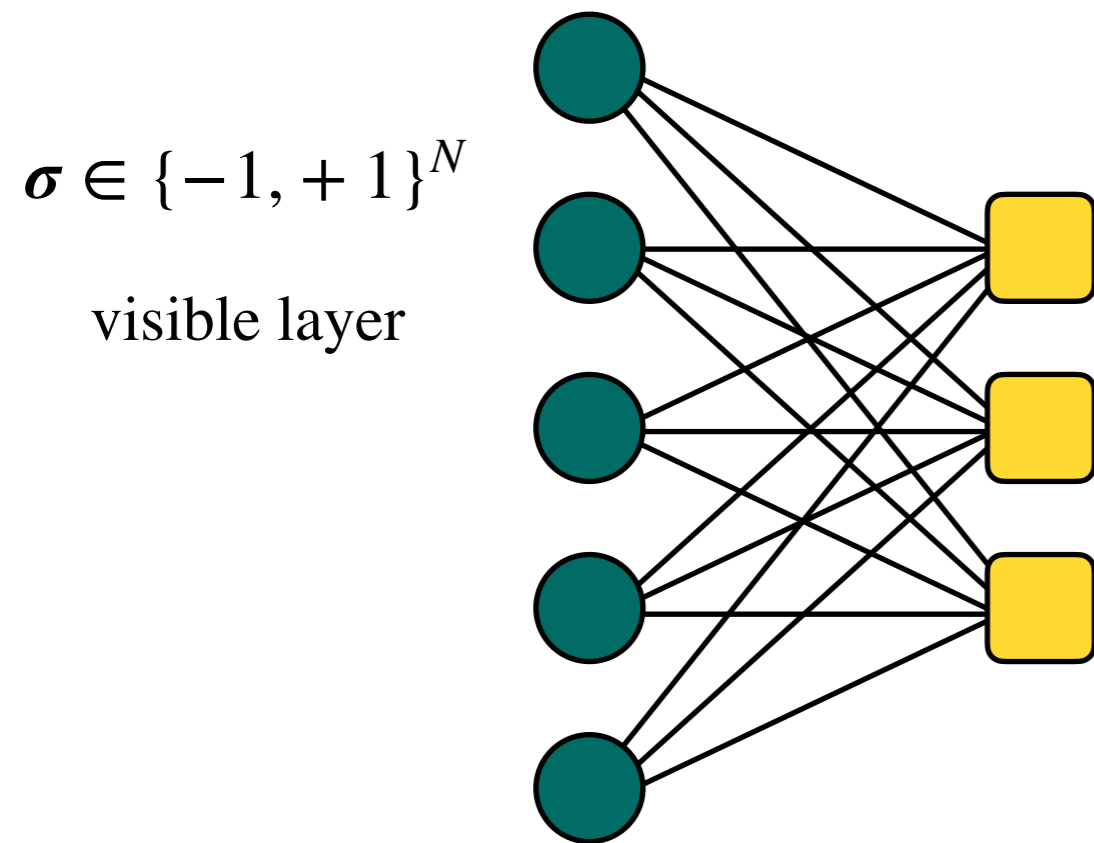
Restricted Boltzmann machines



Restricted Boltzmann machines



Restricted Boltzmann machines



$z \in \mathbb{R}^K$, Gaussian prior

hidden layer

Two layers composed of $N + K$ neurons

$W \in \mathbb{R}^{N \times K}$ interaction matrix (symmetric, zero eye)

Boltzmann-Gibbs measure $\mathcal{P}_W(\sigma, z) \propto e^{\beta \sum_{i,\mu} \sigma_i W_{i,\mu} z_\mu} e^{-\frac{\beta z_\mu^2}{2}}$

Stat-mech: bipartite spin-glass $\mathcal{H}_{N,K,W}^{(\text{RBM})}(\sigma, z) = - \sum_{i,\mu}^{N,K} W_{i,\mu} \sigma_i z_\mu$

Task: classification

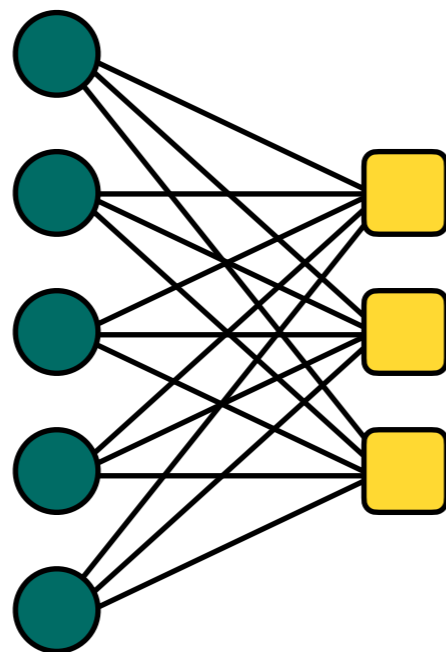
$$\mathcal{S} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(M)}, y^{(M)})\}$$

$$(\mathbf{x}^{(i)}, y^{(i)}) \sim_{iid} \mathcal{Q}(\mathbf{x}, y) \quad \text{unknown, target distr.}$$

$$\hat{y} \sim \mathcal{Q}(y | \mathbf{x})$$






$$\sigma \in \{-1, +1\}^N, \quad N = \text{\#pixels}$$

encodes for \mathbf{x}



$$\mathbf{z} \in \mathbb{R}^K, \quad K = \text{\#classes}$$

encodes for y

	$(\mathbf{x}^{(1)}, 1)$
	$(\mathbf{x}^{(2)}, 3)$
	$(\mathbf{x}^{(3)}, 2)$
	$(\mathbf{x}^{(4)}, 2)$
	$(\mathbf{x}^{(5)}, 1)$
\vdots	

Ideally: Look for a configuration \hat{W} of the matrix W s.t. $\mathcal{P}_{\hat{W}}(\sigma, \mathbf{z}) = \mathcal{Q}(\sigma, \mathbf{z})$, namely s.t.

$$\hat{W} = \min_{W \in \mathbb{R}^{N \times K}} D_{\text{KL}}(\mathcal{Q} \| \mathcal{P})$$

Task: classification

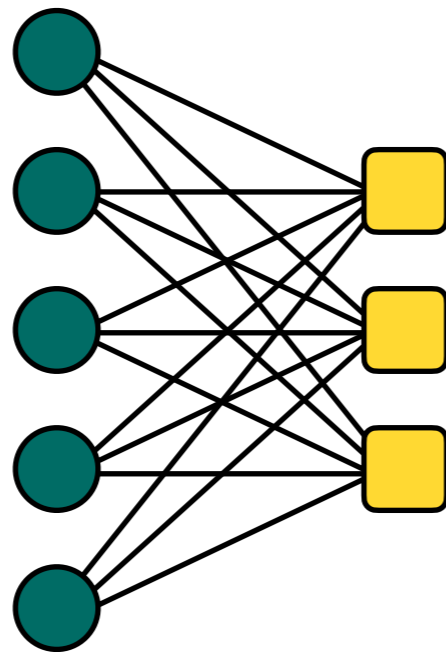
$$\mathcal{S} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(M)}, y^{(M)})\}$$

$$(\mathbf{x}^{(i)}, y^{(i)}) \sim_{iid} \mathcal{Q}(\mathbf{x}, y) \quad \text{unknown, target distr.}$$

$$\hat{y} \sim \mathcal{Q}(y | \mathbf{x})$$






$$\sigma \in \{-1, +1\}^N, \quad N = \text{\#pixels}$$

encodes for \mathbf{x}



$$\mathbf{z} \in \mathbb{R}^K, \quad K = \text{\#classes}$$

encodes for y

	$(\mathbf{x}^{(1)}, 1)$
	$(\mathbf{x}^{(2)}, 3)$
	$(\mathbf{x}^{(3)}, 2)$
	$(\mathbf{x}^{(4)}, 2)$
	$(\mathbf{x}^{(5)}, 1)$
\vdots	

Ideally: Look for a configuration $\hat{\mathbf{W}}$ of the matrix \mathbf{W} s.t. $\mathcal{P}_{\hat{\mathbf{W}}}(\sigma, \mathbf{z}) = \mathcal{Q}(\sigma, \mathbf{z})$, namely s.t.

$$\hat{\mathbf{W}} = \min_{\mathbf{W} \in \mathbb{R}^{N \times K}} D_{\text{KL}}(\mathcal{Q} \| \mathcal{P})$$

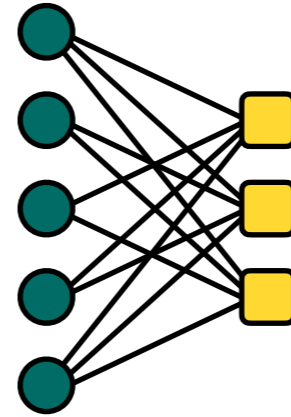
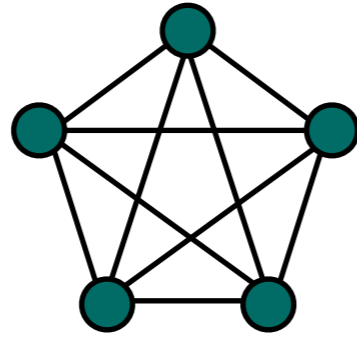
In practice: Look for a configuration $\hat{\mathbf{W}}$ of the matrix \mathbf{W} s.t. $\mathcal{P}_{\hat{\mathbf{W}}}(\sigma, \mathbf{z}) \approx \tilde{\mathcal{Q}}(\sigma, \mathbf{z})$ by applying a gradient

$$\text{descent } W_{i,\mu}^{n+1} = W_{i,\mu}^n - \epsilon \frac{dD_{\text{KL}}(\tilde{\mathcal{Q}} \| \mathcal{P})}{dW_{i,\mu}}$$

Hopfield network / Restricted Boltzmann machines equivalence

$$\sigma \in \{-1, +1\}^N$$

$$\xi \in \{-1, +1\}^{N \times K}$$



$$\sigma \in \{-1, +1\}^N$$

$$z \sim \mathcal{N}(0, \mathbf{I}\beta^{-1})$$

$$W = \xi$$

Hubbard-Stratonovich
transformation



$$\mathcal{Z}_{\beta, N, K}^{(\text{HN})}(\xi) = \sum_{\sigma} e^{\frac{\beta}{2N} \sum_{i, j, \mu} \xi_i^{\mu} \xi_j^{\mu} \sigma_i \sigma_j} = \sum_{\sigma} \int \prod_{\mu} dz_{\mu} e^{-\frac{\beta z_{\mu}^2}{2}} e^{\frac{\beta}{\sqrt{N}} \sum_{i, \mu} \sigma_i \xi_i^{\mu} z_{\mu}} = \mathcal{Z}_{\beta, N, K}^{(\text{RBM})}(\xi)$$



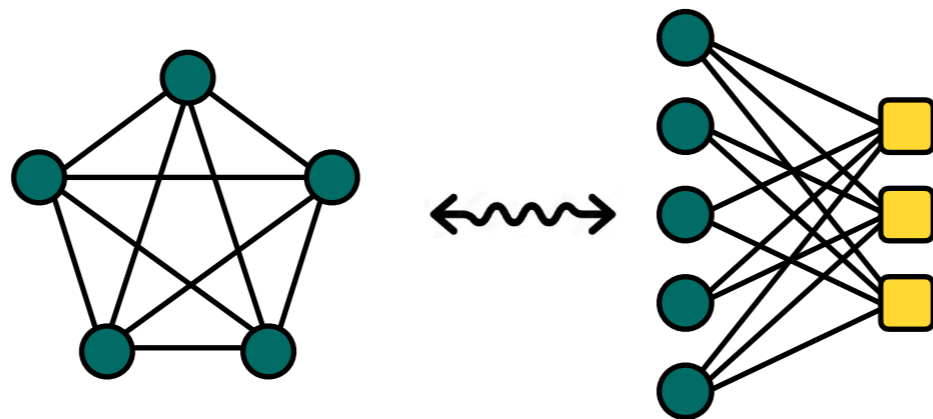
Gaussian
integration

$$\mathcal{F} := -T \log \mathcal{Z} \dots$$

Hopfield network / Restricted Boltzmann machines equivalence

$$\sigma \in \{-1, +1\}^N$$

$$\xi \in \{-1, +1\}^{N \times K}$$



$$\sigma \in \{-1, +1\}^N$$

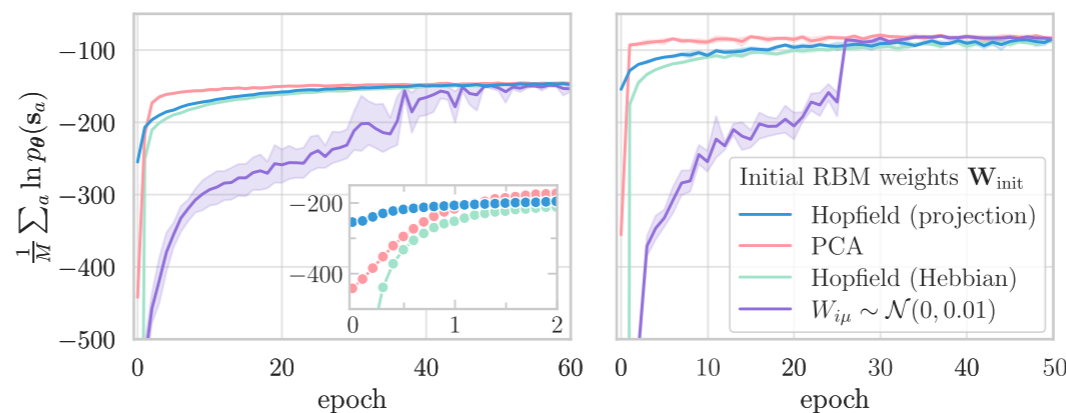
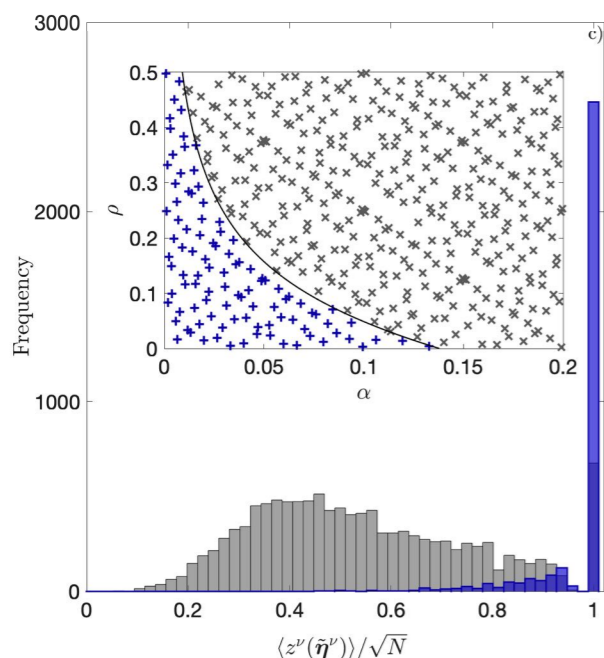
$$z \sim \mathcal{N}(0, \mathbf{I}\beta^{-1})$$

$$\mathbf{W} = \xi$$

Hubbard-Stratonovich transformation

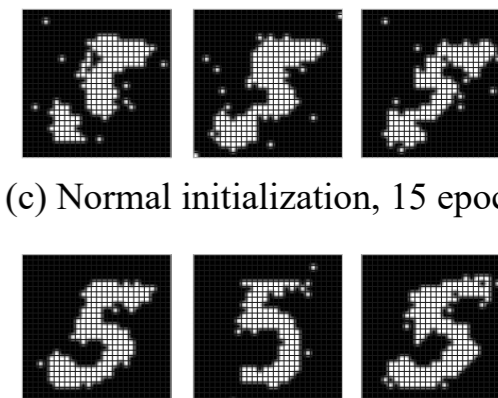
$$\mathcal{L}_{\beta, N, K}^{(\text{HN})}(\xi) = \sum_{\sigma} e^{\frac{\beta}{2N} \sum_{i, j, \mu} \xi_i^{\mu} \xi_j^{\mu} \sigma_i \sigma_j} = \sum_{\sigma} \int \prod_{\mu} dz_{\mu} e^{-\frac{\beta z_{\mu}^2}{2}} e^{\frac{\beta}{\sqrt{N}} \sum_{i, \mu} \sigma_i \xi_i^{\mu} z_{\mu}} = \mathcal{L}_{\beta, N, K}^{(\text{RBM})}(\xi)$$

Gaussian integration



(a) 10 hidden units

(b) 50 hidden units



(c) Normal initialization, 15 epochs

(d) Hopfield initialization, 15 epochs

Dataset and architecture complexities interplay

MNIST and Fashion-MNIST Datasets

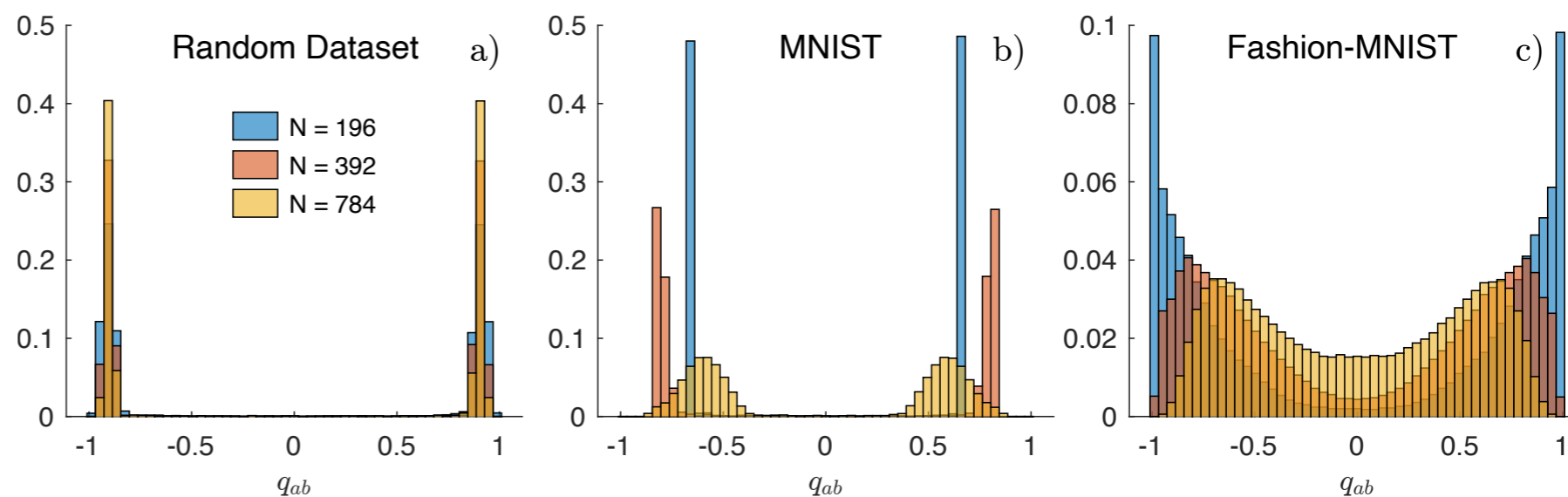
$K = 10$ classes

$M = 6000$ examples each



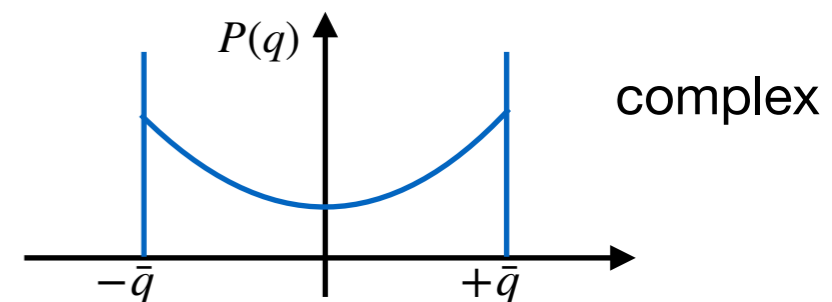
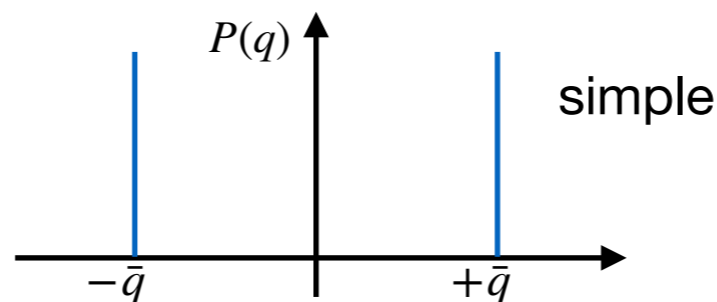
example overlaps

$$\tilde{q}_{ab} = \frac{1}{N} \sum_{i=1}^N \eta_i^{\mu,a} \eta_i^{\mu,b}$$



replica overlaps

$$q_{ab} = \frac{1}{N} \sum_{i=1}^N \sigma_i^{(a)} \sigma_i^{(b)}$$



Dataset and architecture complexities interplay

MNIST and Fashion-MNIST Datasets

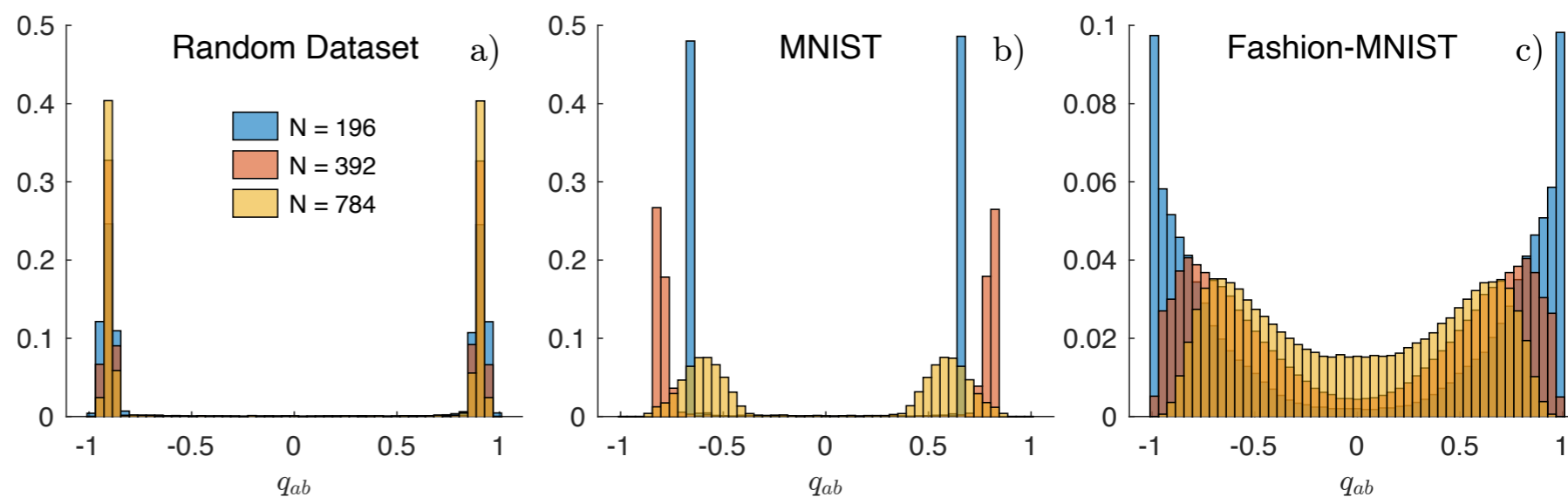
$K = 10$ classes

$M = 6000$ examples each



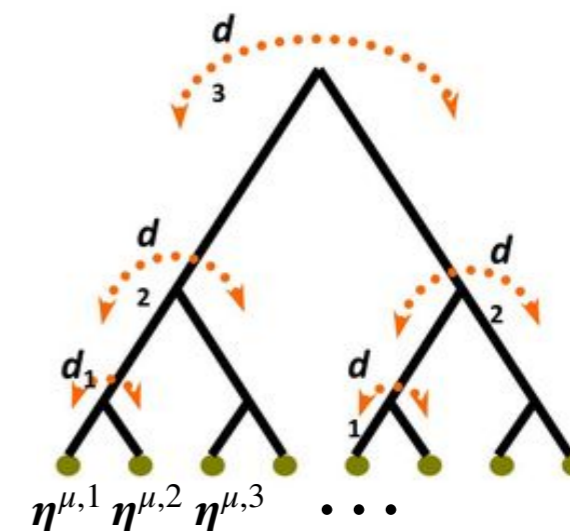
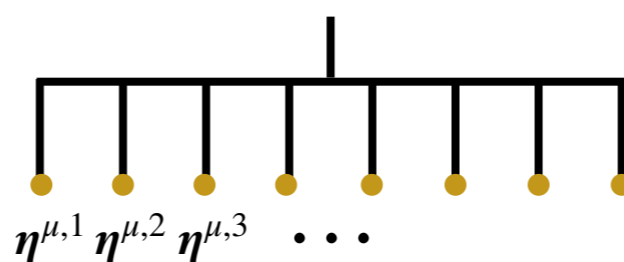
example overlaps

$$\tilde{q}_{ab} = \frac{1}{N} \sum_{i=1}^N \eta_i^{\mu,a} \eta_i^{\mu,b}$$



replica overlaps

$$q_{ab} = \frac{1}{N} \sum_{i=1}^N \sigma_i^{(a)} \sigma_i^{(b)}$$



Dataset and architecture complexities interplay

MNIST and Fashion-MNIST Datasets

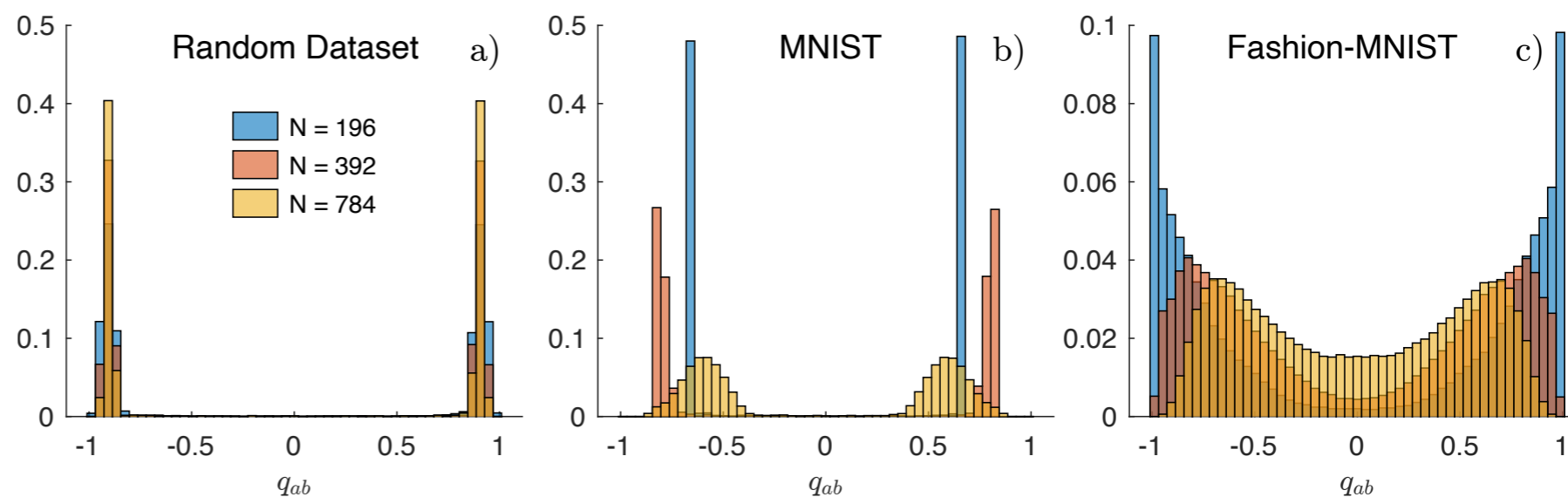
$K = 10$ classes

$M = 6000$ examples each



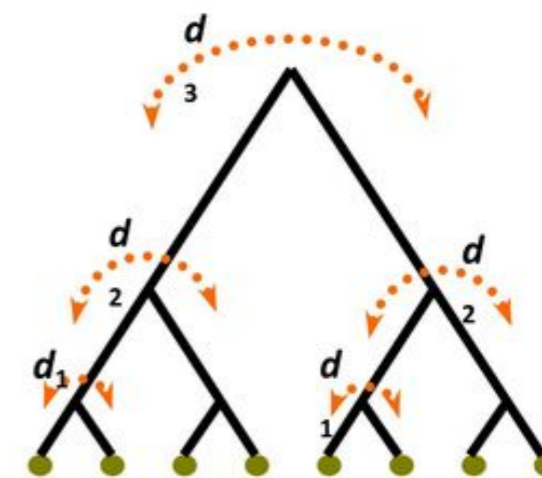
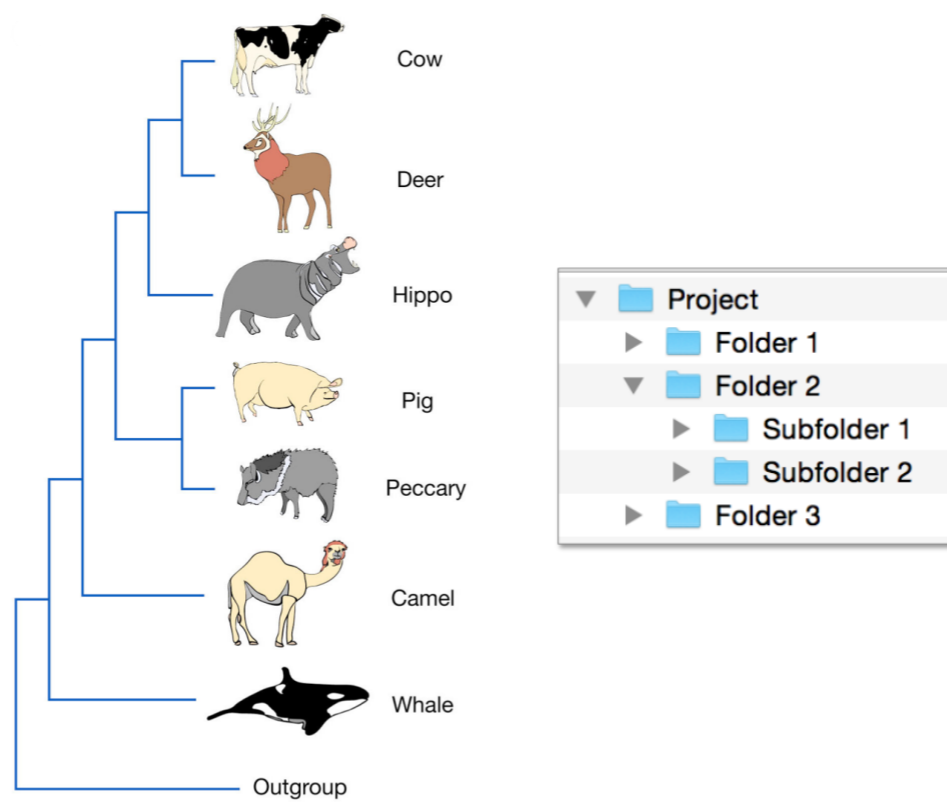
example overlaps

$$\tilde{q}_{ab} = \frac{1}{N} \sum_{i=1}^N \eta_i^{\mu,a} \eta_i^{\mu,b}$$



replica overlaps

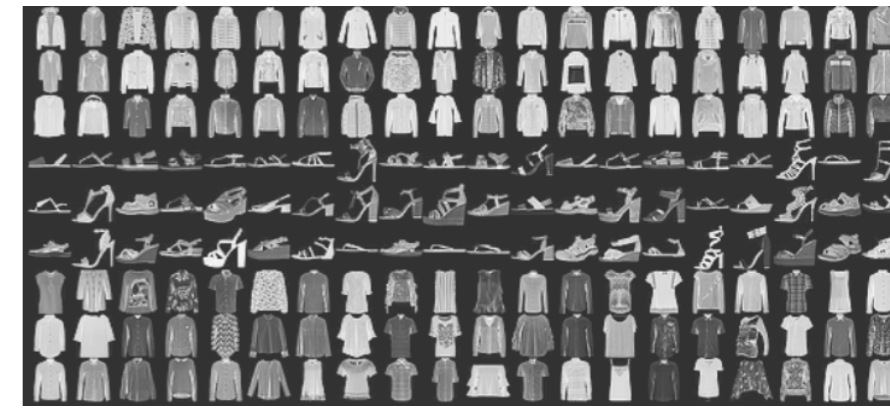
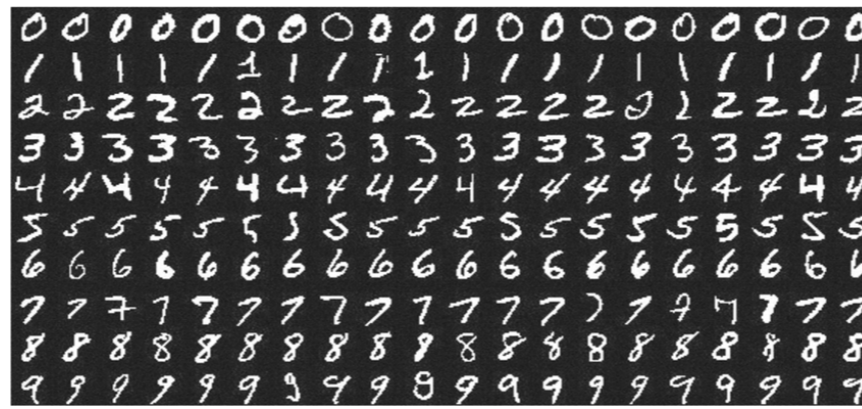
$$q_{ab} = \frac{1}{N} \sum_{i=1}^N \sigma_i^{(a)} \sigma_i^{(b)}$$



MNIST and Fashion-MNIST Datasets

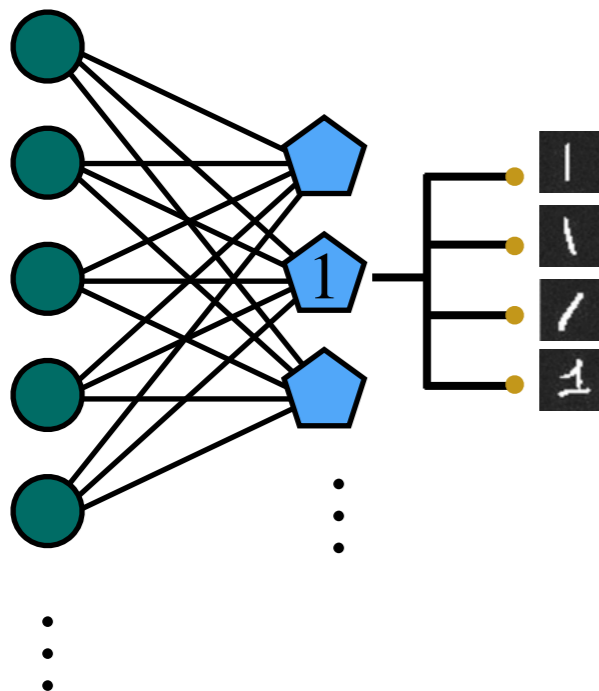
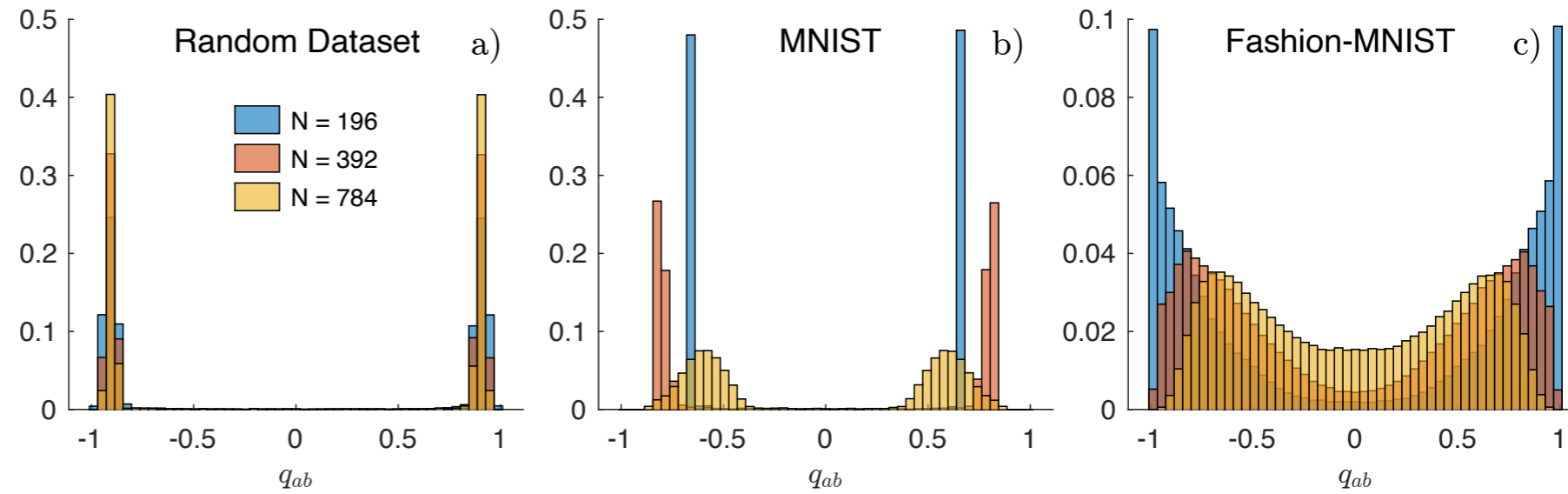
$K = 10$ classes

$M = 6000$ examples each



example overlaps

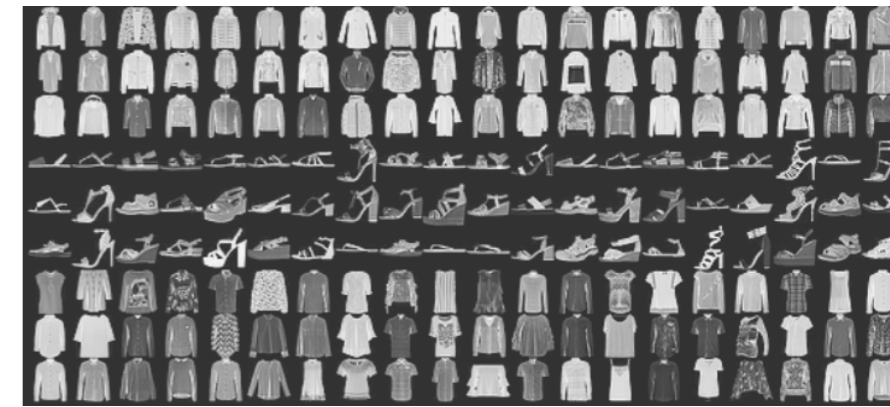
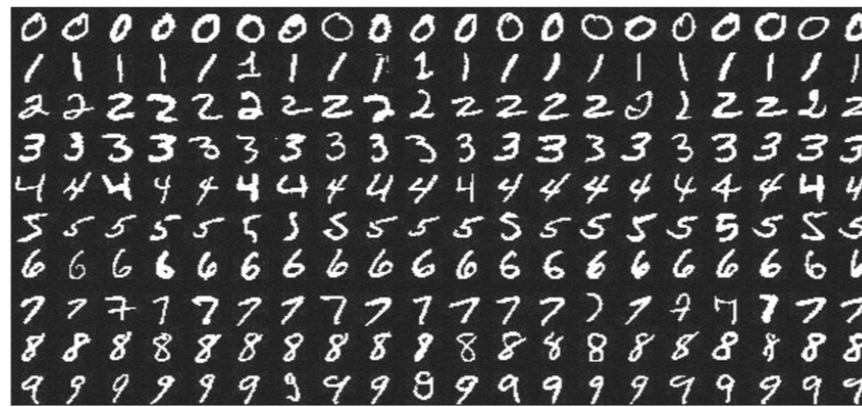
$$\tilde{q}_{ab} = \sum_{i=1}^N \xi_i^{(a)} \xi_i^{(b)}$$



MNIST and Fashion-MNIST Datasets

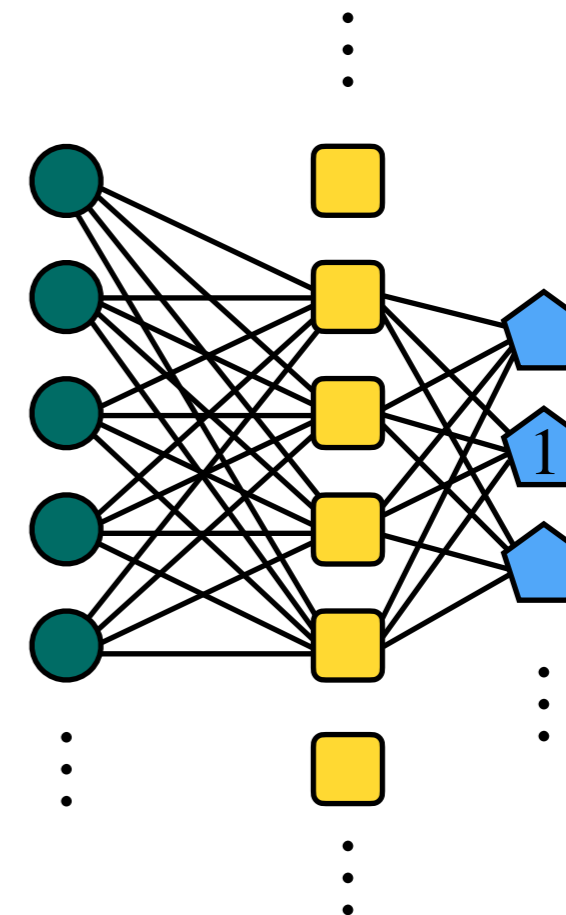
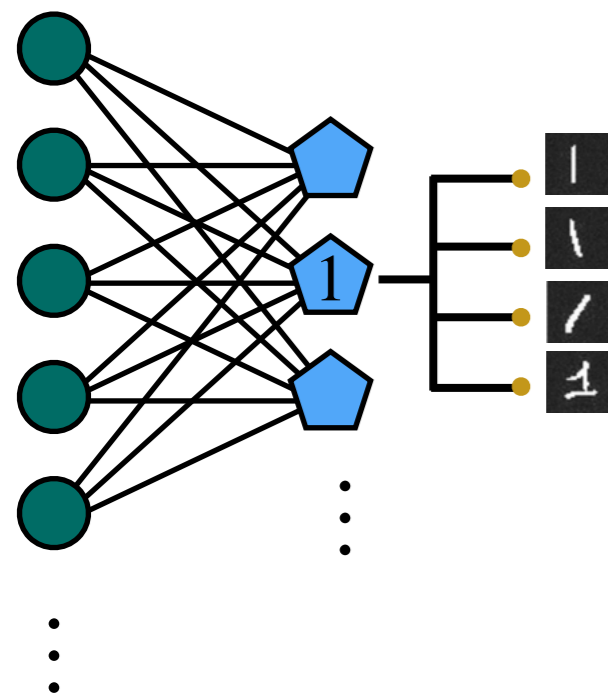
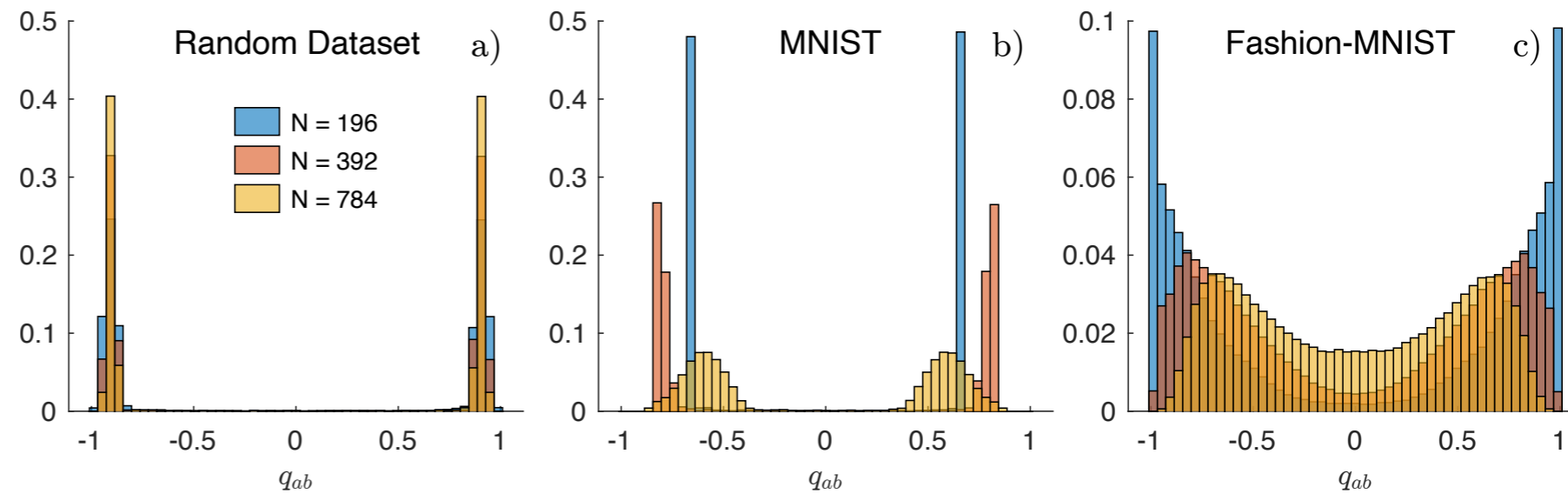
$K = 10$ classes

$M = 6000$ examples each



example overlaps

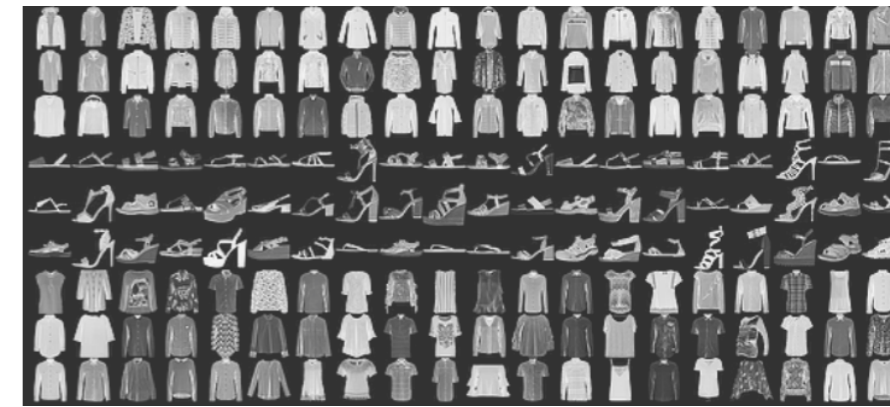
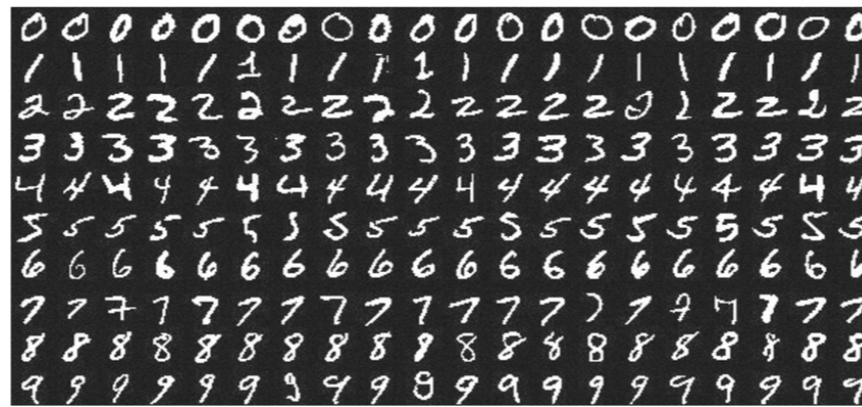
$$\tilde{q}_{ab} = \sum_{i=1}^N \xi_i^{(a)} \xi_i^{(b)}$$



MNIST and Fashion-MNIST Datasets

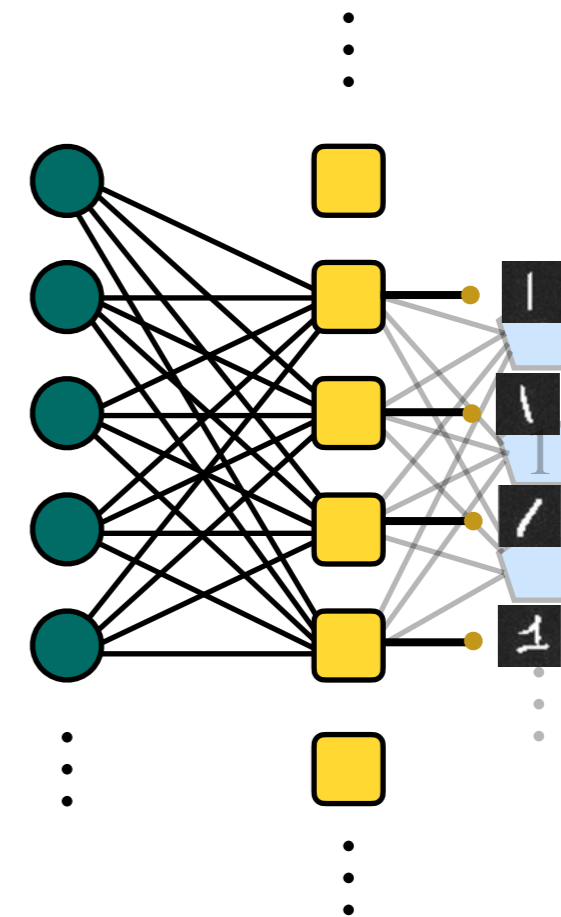
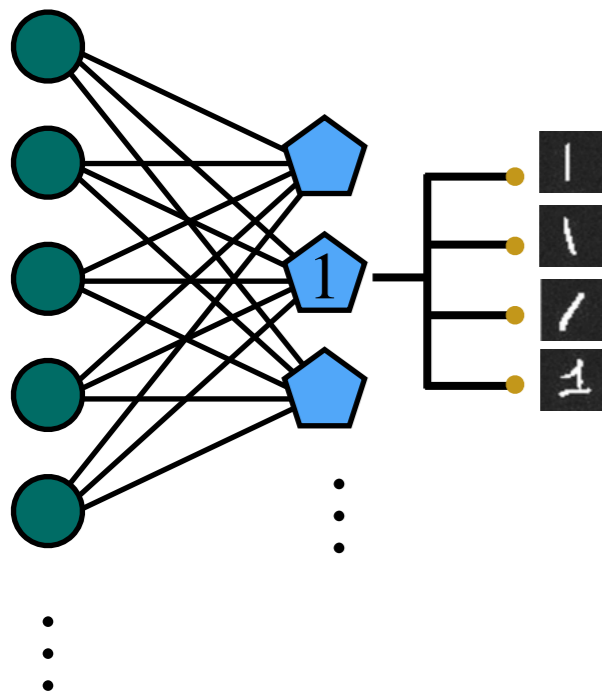
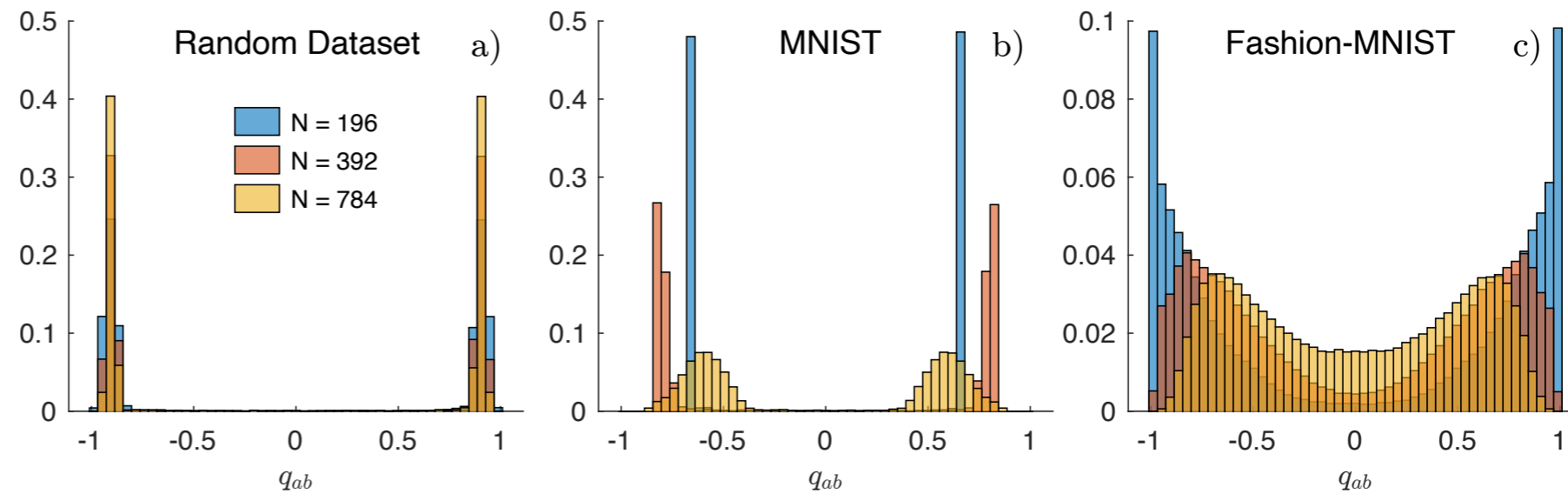
$K = 10$ classes

$M = 6000$ examples each



example overlaps

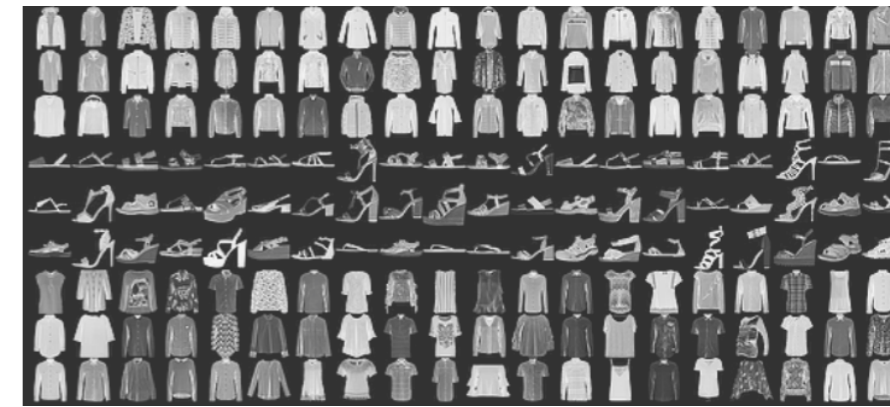
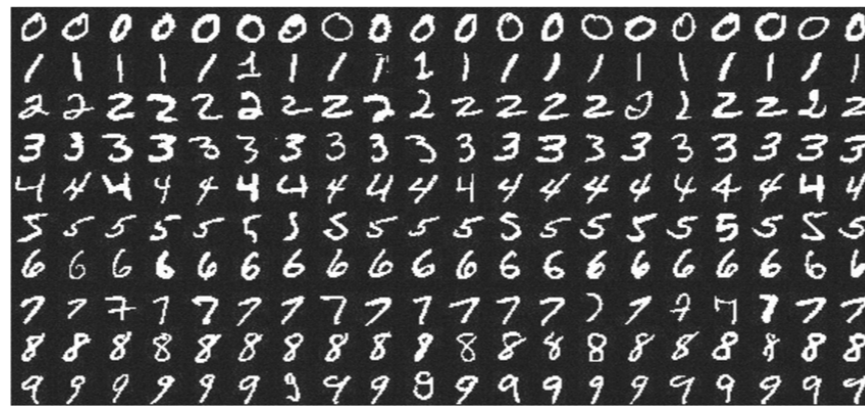
$$\tilde{Q}_{ab} = \sum_{i=1}^N \xi_i^{(a)} \xi_i^{(b)}$$



MNIST and Fashion-MNIST Datasets

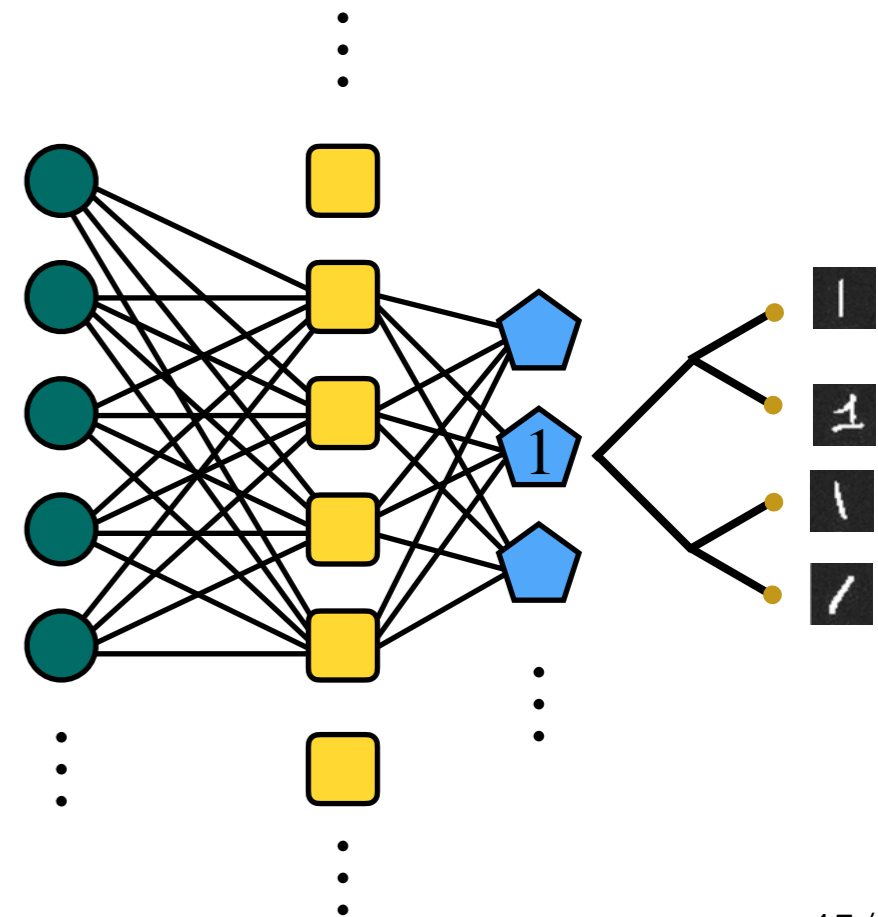
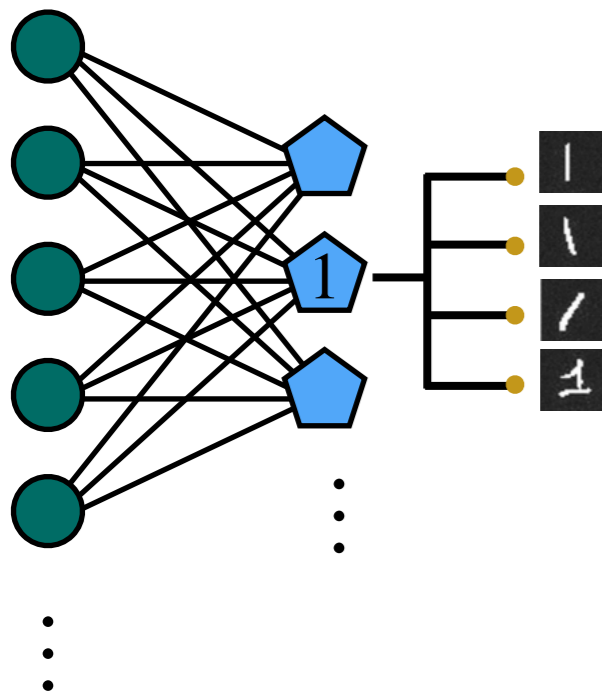
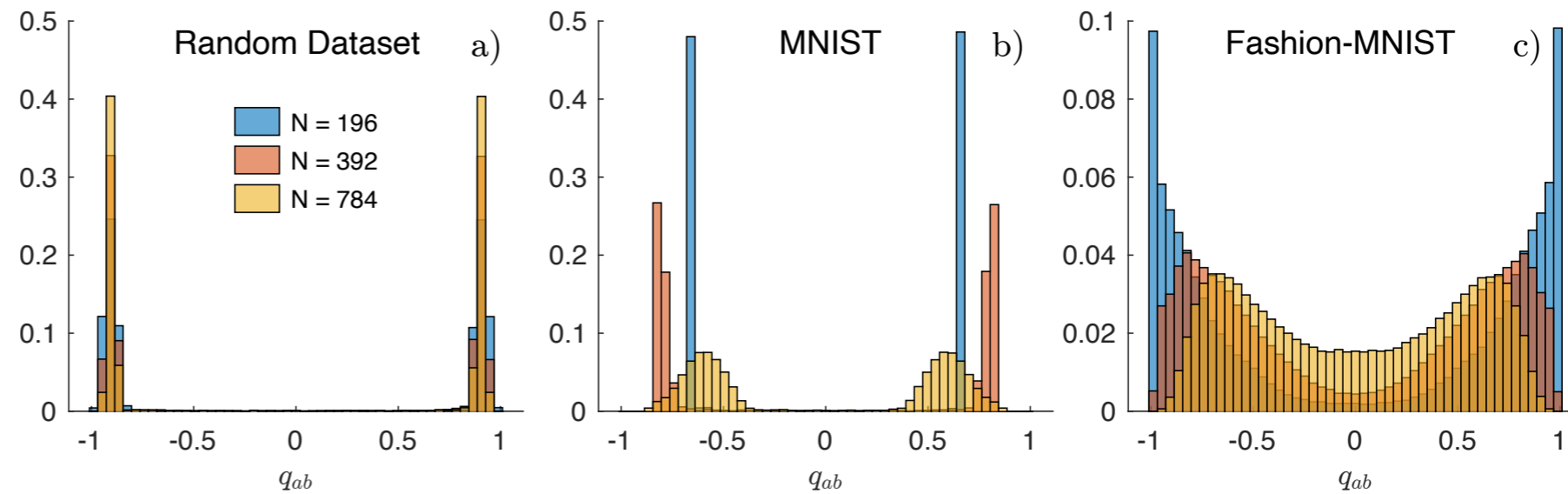
$K = 10$ classes

$M = 6000$ examples each



example overlaps

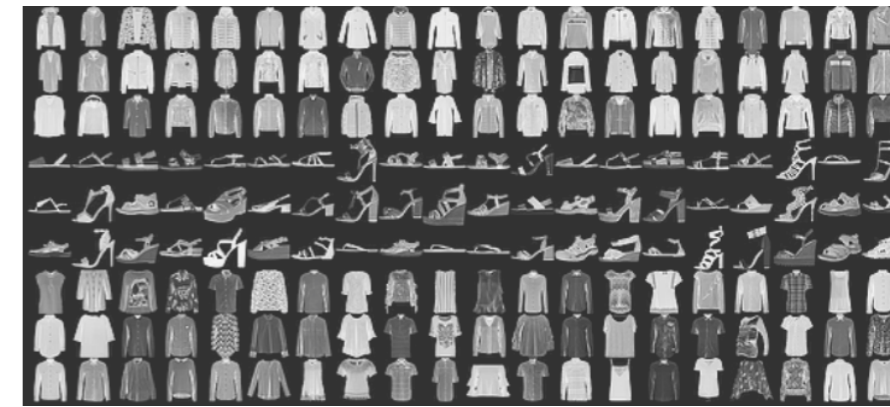
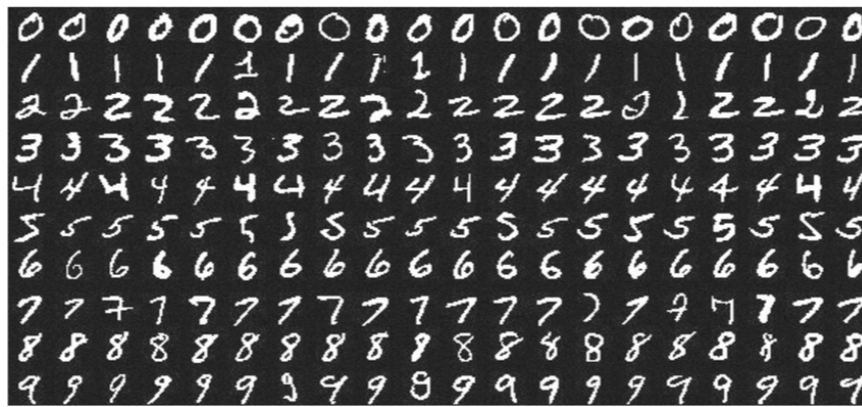
$$\tilde{q}_{ab} = \sum_{i=1}^N \xi_i^{(a)} \xi_i^{(b)}$$



MNIST and Fashion-MNIST Datasets

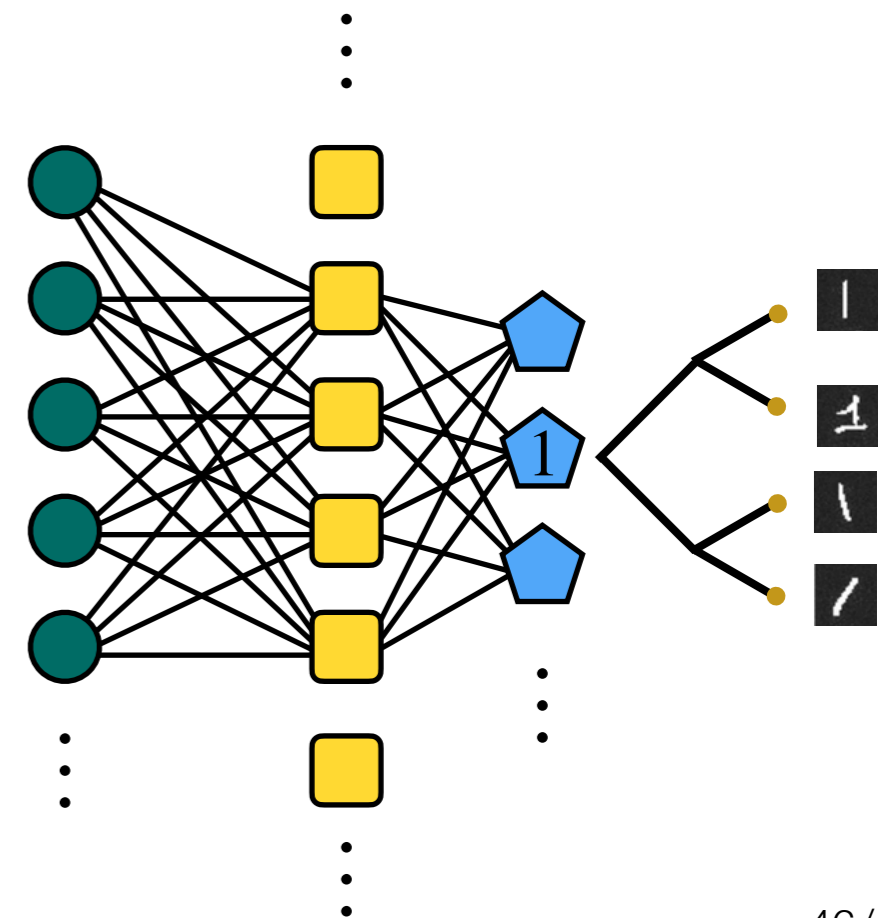
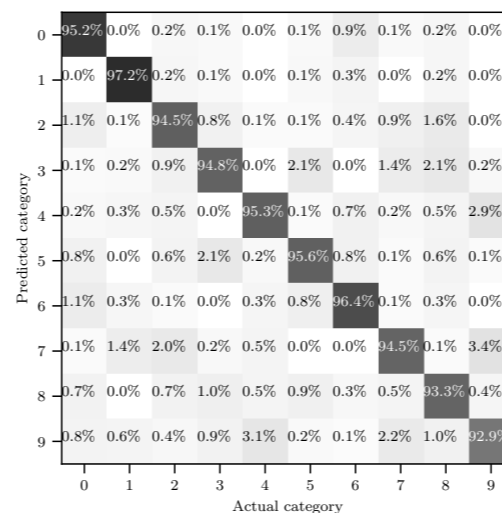
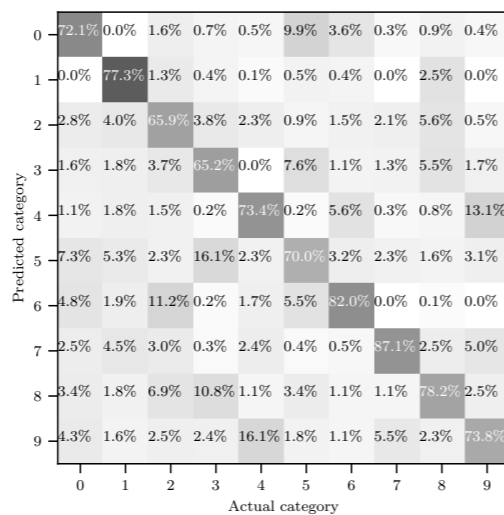
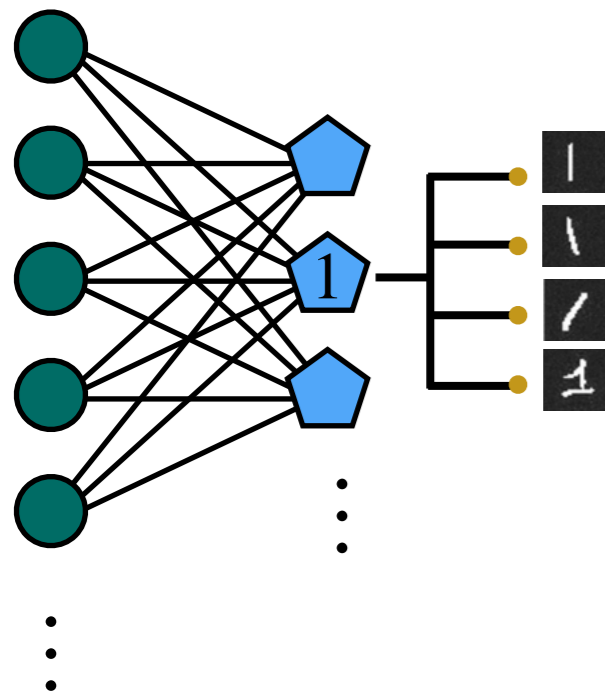
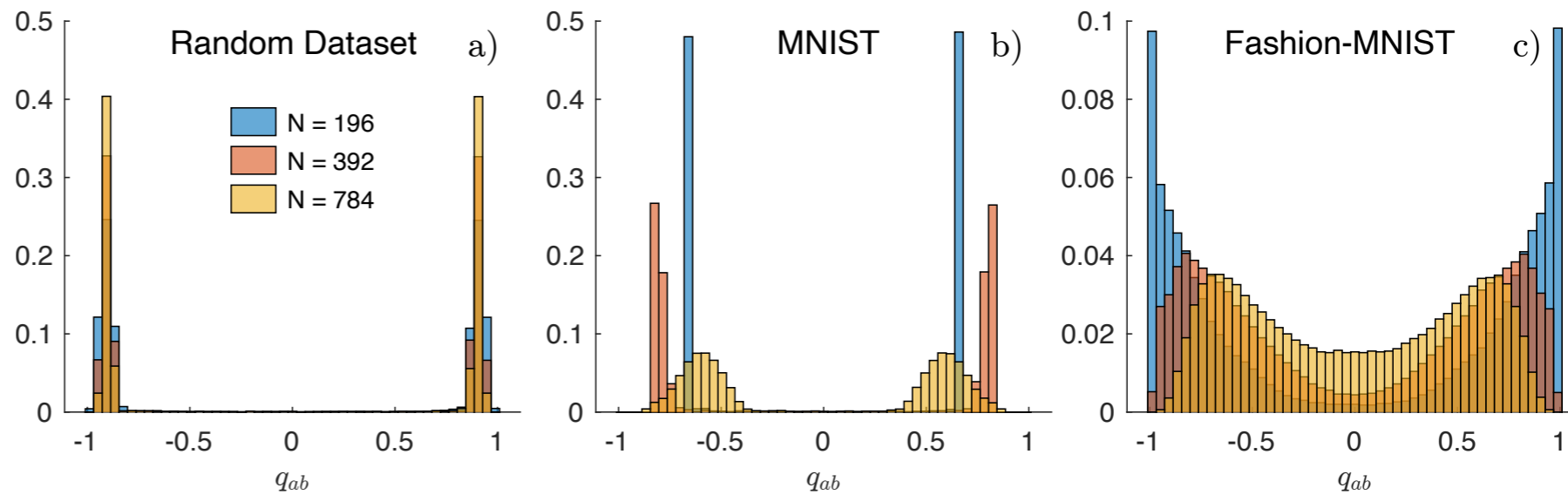
$K = 10$ classes

$M = 6000$ examples each



example overlaps

$$\tilde{q}_{ab} = \sum_{i=1}^N \xi_i^{(a)} \xi_i^{(b)}$$



THANKS

