#rbm_tutorial

@Sehmimul Hoque

Ejaaz Merali

@emerali
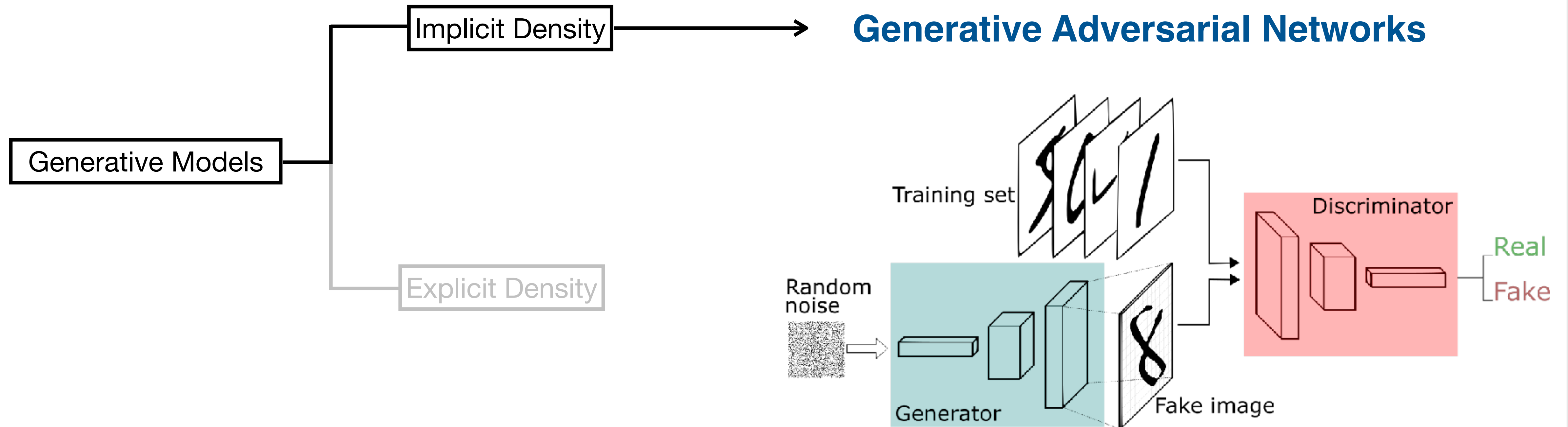
QUANTUM INTELLIGENCE LAB

UNIVERSITY OF WATERLOO

# Taxonomy of generative models



Generative Models
- Implicit Density
- Explicit Density
  - Approximate Density
  - Tractable Density

I. Goodfellow arXiv:1701.00160
M. Albergo https://machine19.github.io

# Taxonomy of generative models



**Generative Models** — **Implicit Density** → **Generative Adversarial Networks**

Explicit Density

# Taxonomy of generative models



**Hopfield Networks**
**Restricted Boltzmann Machines**
**Variational Autoencoders**

Generative Models

Implicit Density

Explicit Density

Approximate Density

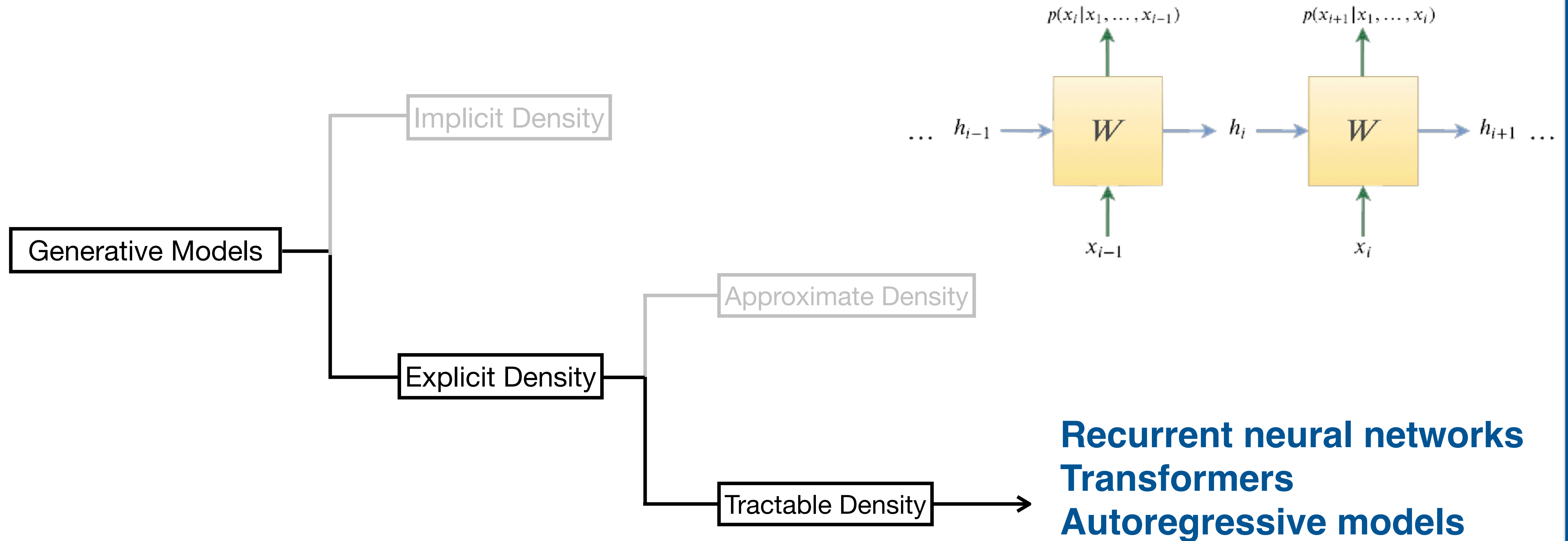Tractable Density

$x_i = 0, 1$

$h_j = 0, 1$

Torlai and RGM, Phys. Rev. B 94, 165134 (2016)
Carleo, Troyer Science 355, 602 (2017)
S. Wetzel, arXiv:1703.02435
RGM, Carleo, Carrasquilla, Cirac, Nature Physics 15, 887 (2019)

# Taxonomy of generative models

$$p(x_i \mid x_1, \ldots, x_{i-1}) \qquad p(x_{i+1} \mid x_1, \ldots, x_i)$$

$$\ldots \; h_{i-1} \longrightarrow \boxed{W} \longrightarrow h_i \longrightarrow \boxed{W} \longrightarrow h_{i+1} \; \ldots$$

$$x_{i-1} \qquad\qquad x_i$$

Generative Models

Implicit Density

Explicit Density

Approximate Density

Tractable Density

**Recurrent neural networks**
**Transformers**
**Autoregressive models**

Carrasquilla, Torlai, RGM, Aolita, Nature Machine Intelligence 1, 155 (2019)
Sharir, Levine, Wies, Carleo, Shashua, Phys. Rev. Lett. 124, 020503 (2020)
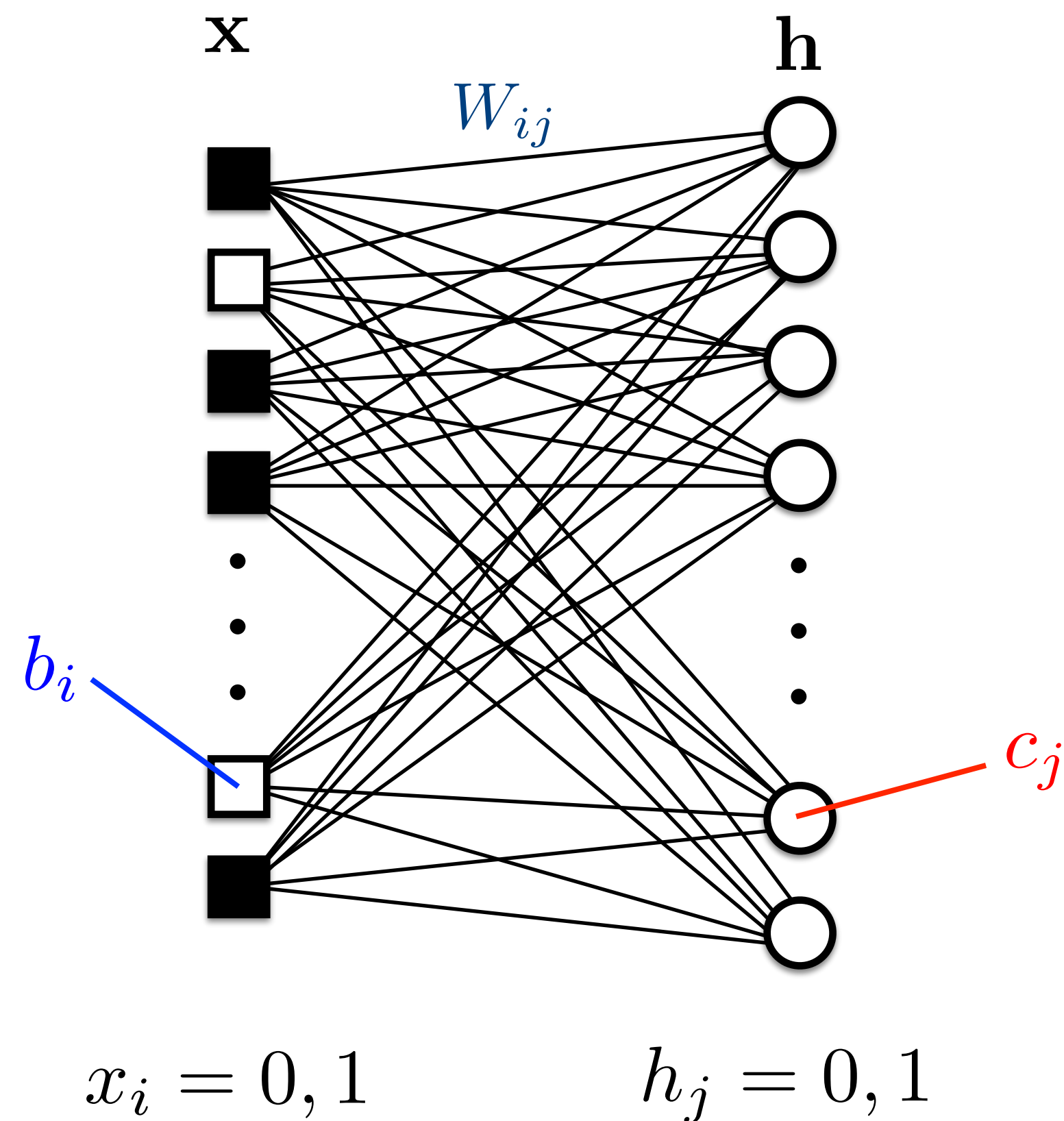Cha, Ginsparg, Wu, Carrasquilla, L. McMahon, Kim, arXiv:2006.12469
Hibat-Allah, Ganahl, Hayward, RGM, and Carrasquilla, Phys. Rev. Research 2, 023358 (2020)

# Restricted Boltzmann Machine

Smolensky, Hinton, Salakhutdinov, Bengio

$N$ visible units    $n_h$ hidden units



$x_i = 0, 1$    $h_j = 0, 1$

Like a Hopfield network, RBMs are "energy-based" models:

$$p_\lambda = \frac{1}{Z_\lambda} e^{-E_\lambda(\mathbf{x}, \mathbf{h})}$$

joint probability distribution

$$E_\lambda(\mathbf{x}, \mathbf{h}) = -\sum_{ij} W_{ij} x_i h_j - \sum_i b_i x_i - \sum_j c_j h_j$$

"Training" means tuning the machine parameters to get the marginal distribution $p_\lambda(\mathbf{x})$ to approximate the (unknown) target distribution
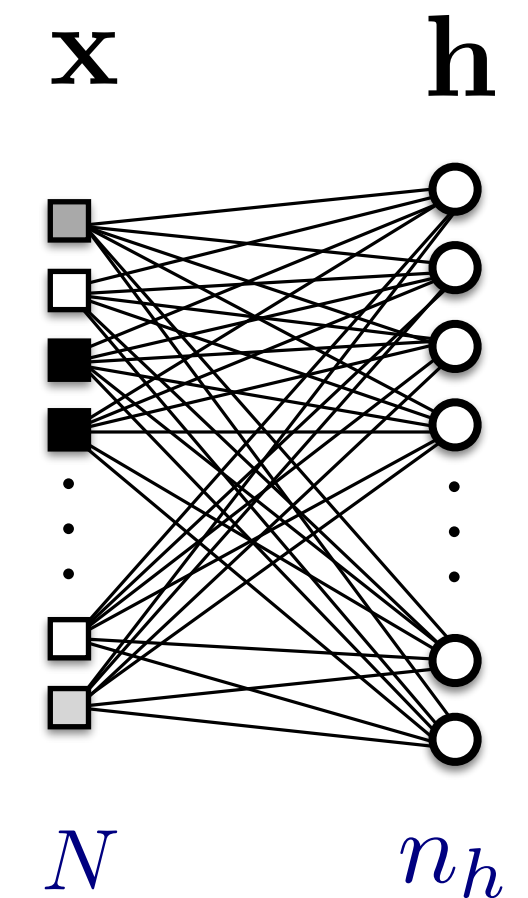
$$\lambda = \{W, b, c\} \qquad p_\lambda(\mathbf{x}) = \sum_{\mathbf{h}} p_\lambda(\mathbf{x}, \mathbf{h})$$

model parameters

# Block Gibbs Sampling

$$\mathbf{x} \qquad \mathbf{h}$$



$$N \qquad n_h$$

RBM: being "restricted" is a special property that allows sampling one layer at a time.

$$x_0 \rightarrow h_0 \rightarrow x_1 \rightarrow h_1 \rightarrow \cdots \rightarrow x_k \rightarrow h_k$$

Each layer is updated with **conditional** probabilities

$$p_\lambda(\mathbf{x}|\mathbf{h}) = \frac{p_\lambda(\mathbf{x}, \mathbf{h})}{p_\lambda(\mathbf{h})} = \prod_i p(x_i|\mathbf{h}) \qquad p_\lambda(\mathbf{h}|\mathbf{x}) = \frac{p_\lambda(\mathbf{x}, \mathbf{h})}{p_\lambda(\mathbf{x})} = \prod_j p(h_j|\mathbf{x})$$
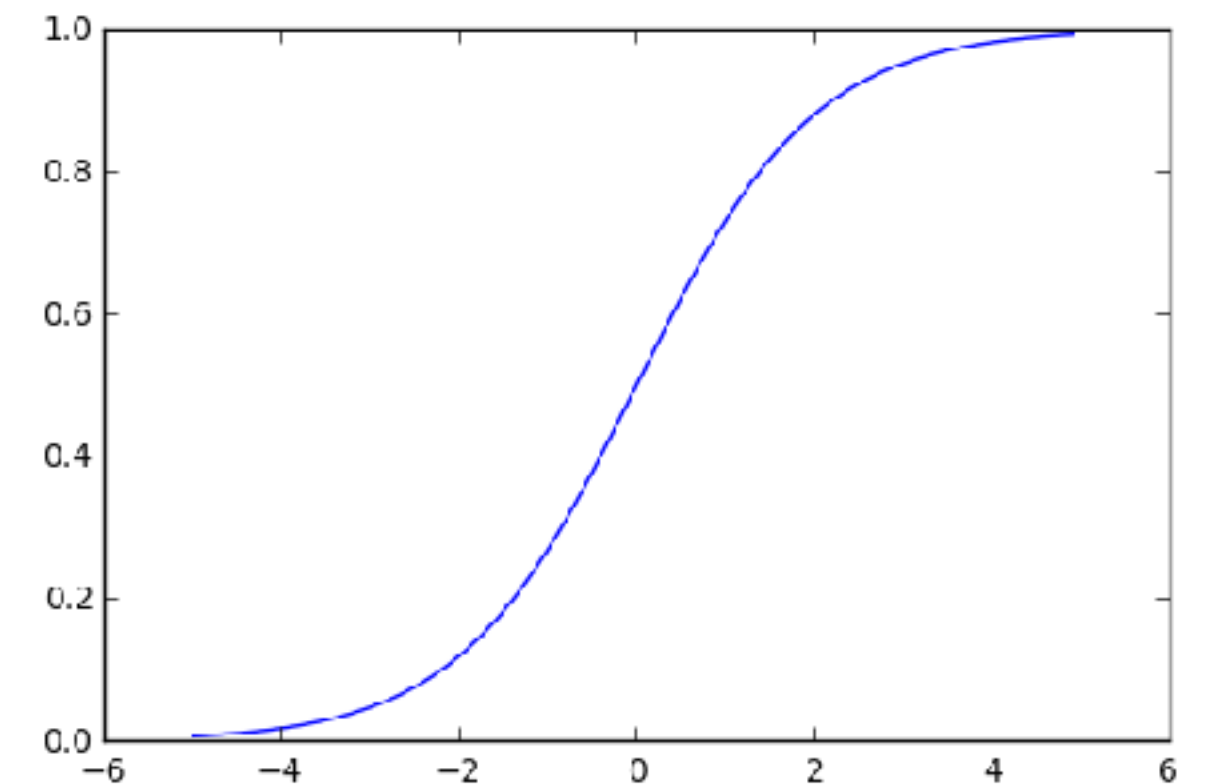
Sigmoid function

$$p(x_i = 1|\mathbf{h}) = \sigma\left(\sum_j W_{ij} h_j + b_i\right)$$

$$p(h_j = 1|\mathbf{x}) = \sigma\left(\sum_i W_{ij} x_i + c_j\right)$$

$\Big\}$ The firing rate of a (stochastic) neuron

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

# Training the RBM

***Training*** means tuning the machine parameters to minimize the difference between the marginal distribution $p_\lambda(\mathbf{x}) = \sum_{\mathbf{h}} p_\lambda(\mathbf{x}, \mathbf{h})$ and the (unknown) physical "target" distribution

Define an optimization problem: minimize the ***Kullback-Leibler divergence***

$$\mathrm{KL}(p||p_\lambda) = \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{p_\lambda(\mathbf{x})} \geq 0$$



$p(\mathbf{x})$     $p_\lambda(\mathbf{x})$

- A non-symmetric measure of the distance between two distributions
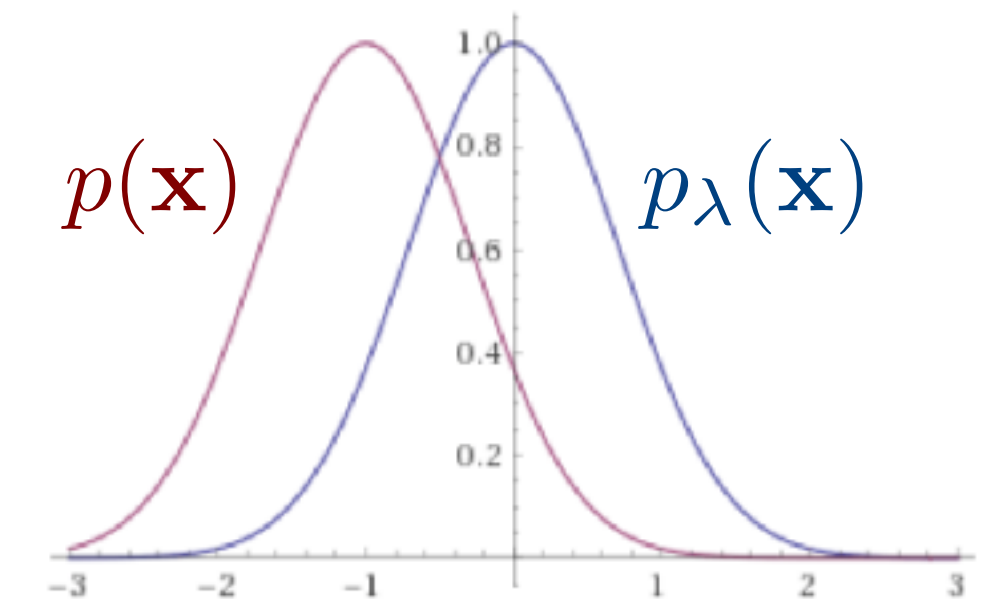- Always positive, and zero iff $p = p_\lambda$

$$\mathrm{KL}(p||p_\lambda) = \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_{\mathbf{x}} p(\mathbf{x}) \log p_\lambda(\mathbf{x})$$

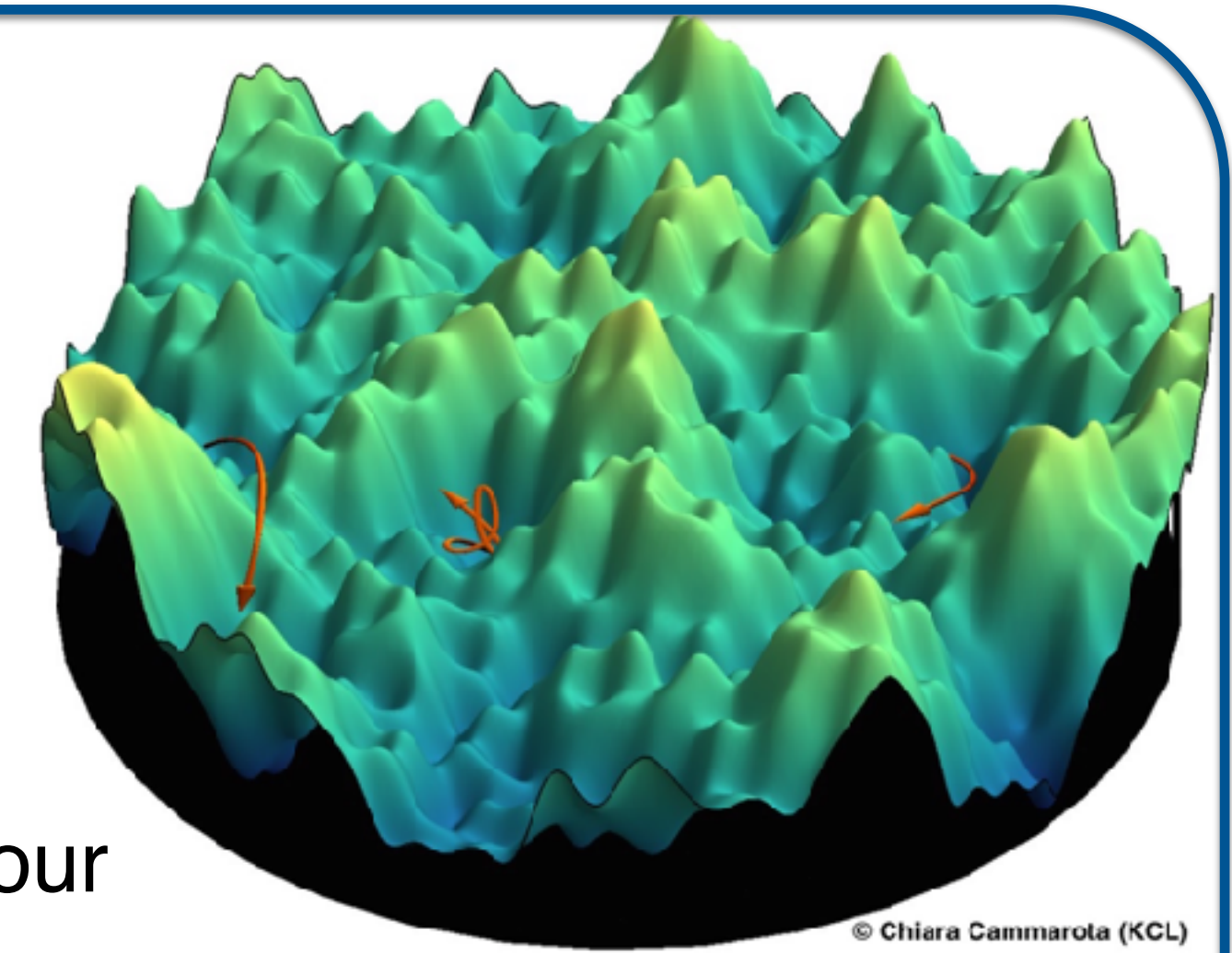the entropy of $p$      depends on the parameters over which to optimize

$$= -\langle \log p_\lambda(\mathbf{x}) \rangle_p \approx -\sum_i \log p_\lambda(\mathbf{x}_i)$$

Equivalent to maximizing the "log-likelihood"    $\mathcal{L} = \langle \log p_\lambda(\mathbf{x}) \rangle_p$

# Stochastic Gradient Descent


© Chiara Cammarota (KCL)

The optimization landscape is thus obtained - minimize using gradient descent

$$\lambda' = \lambda - \eta \nabla \mathcal{L} \qquad \lambda = \{W, b, c\}$$

The full gradient is too costly to calculate. Instead sample some number $m$ of your dataset, and perform *stochastic* gradient descent

$$\frac{1}{m} \sum_{j=1}^{m} \nabla \mathcal{L}(\mathbf{x}_j) \approx \nabla \mathcal{L} \qquad \text{"mini-batch" size = } m$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -\sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{x}) \frac{\partial E}{\partial \lambda} + \sum_{\mathbf{x},\mathbf{h}} p(\mathbf{x}, \mathbf{h}) \frac{\partial E}{\partial \lambda}$$

- The first term is computationally easy to calculate.

- The second term is hard. Requires a MCMC to generate samples from the station distribution of the machine. In practice a short chain of $k$ steps is run
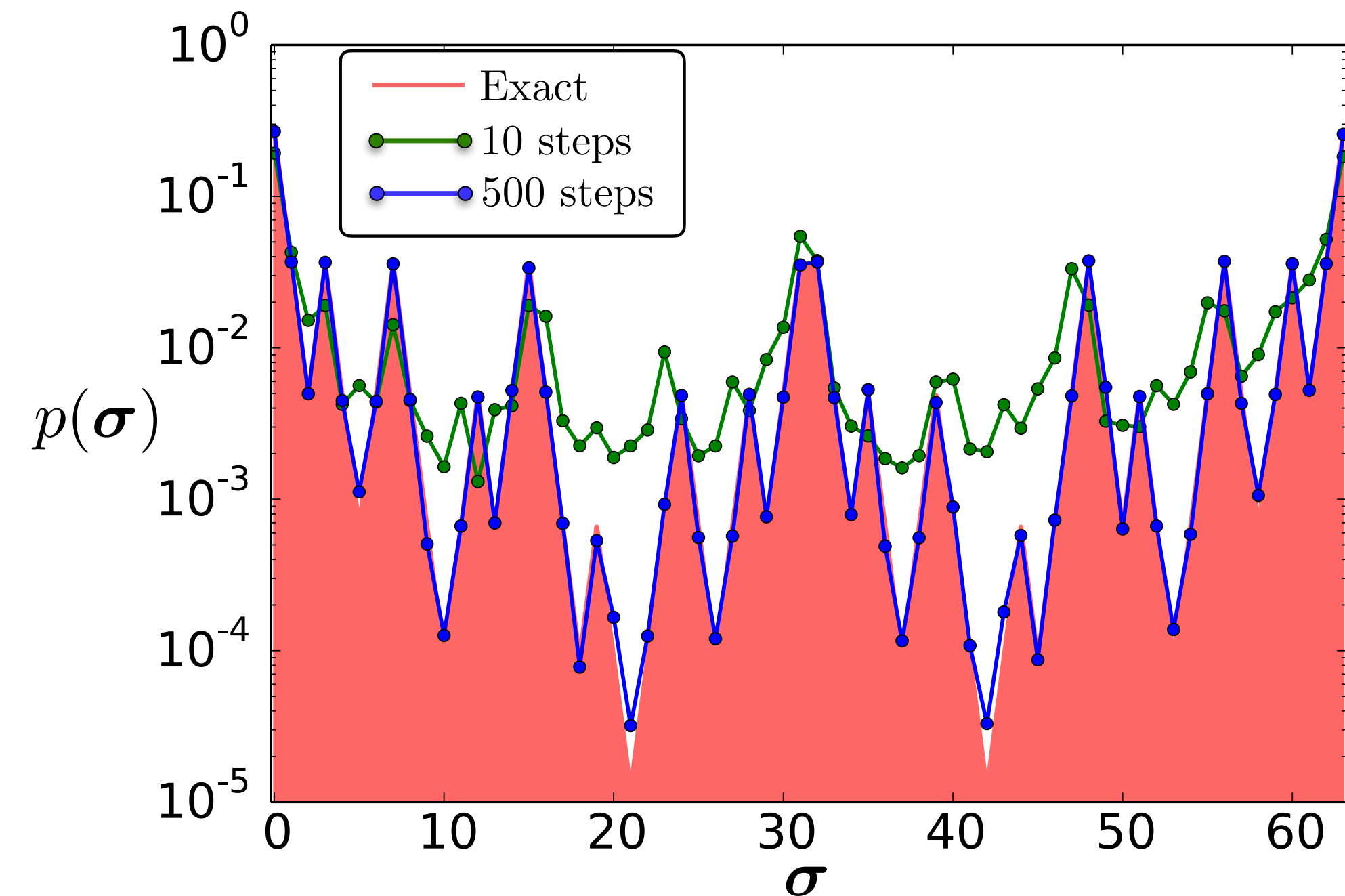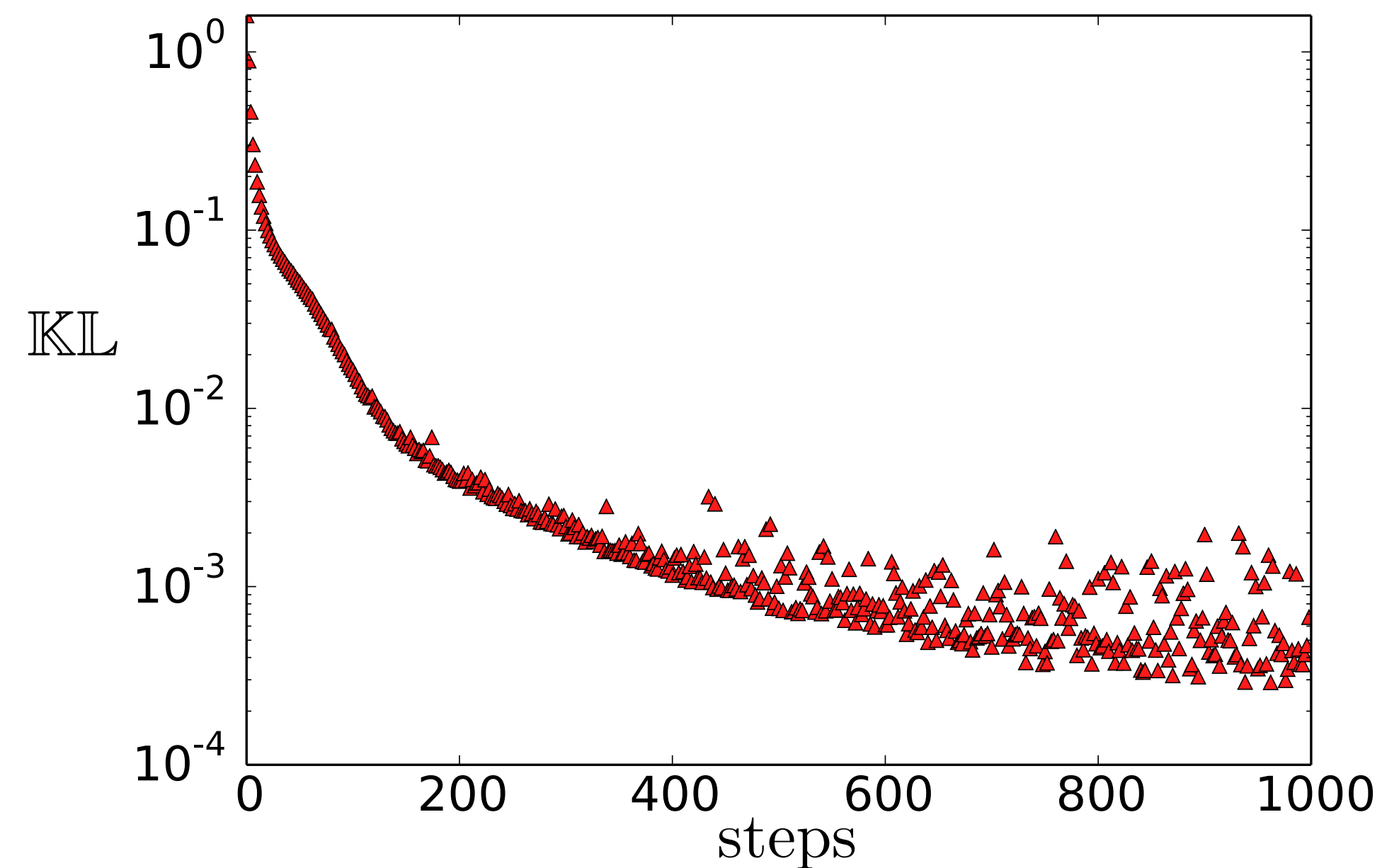
**"Contrastive Divergence"**

$$\mathrm{CD}_k$$

# Stochastic Gradient Descent

Example: 1D Ising model with 6 spins, trained using $CD_5$



KL not possible to calculate in the general case for larger $N$.

Other metrics, like physical observables, could be used to validate the training…

# Sampling the RBM

After the training is complete, and the weights and biases are stationary, one can use Block Gibbs sampling to produce new configurations

$$x_0 \to h_0 \to x_1 \to h_1 \to \cdots \to x_k \to h_k$$

$$\langle \mathcal{O} \rangle_{\text{joint}} = \frac{1}{Z} \sum_{x,h} \mathcal{O} \cdot p_\lambda(x,h)$$

Restrict the observable to the visible layer only $\quad \mathcal{O} = \mathcal{O}_x$

$$\langle \mathcal{O} \rangle_{\text{joint}} = \sum_x \mathcal{O}_x \boxed{\sum_h p_\lambda(x,h)} \approx \langle \mathcal{O} \rangle_{\text{physical}}$$

$$p_\lambda(x) \to p(x)$$

Can ask, how well do physical estimators calculated in this way match the exact values? How does this depend on the number of parameters in the machine for a given system size,
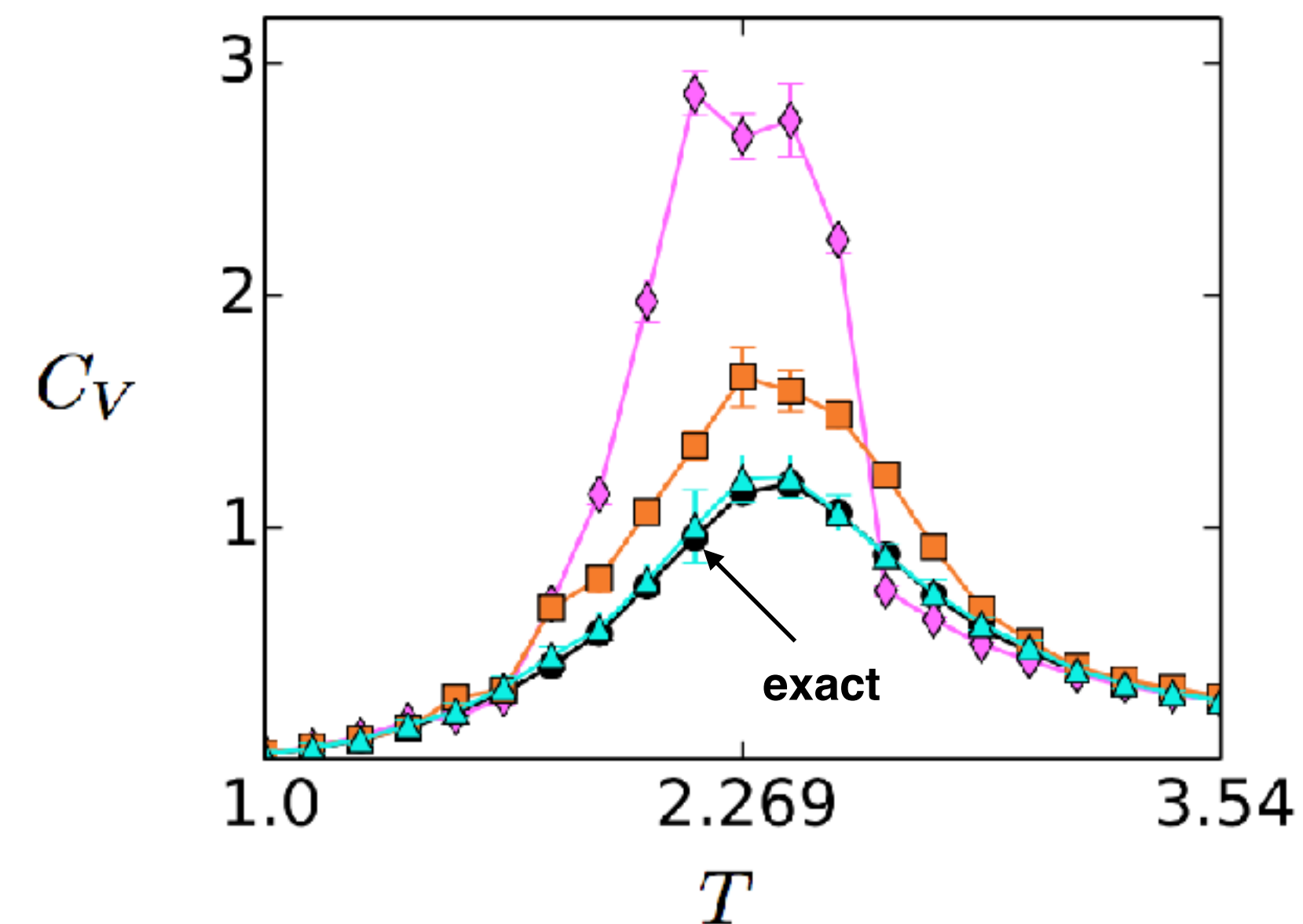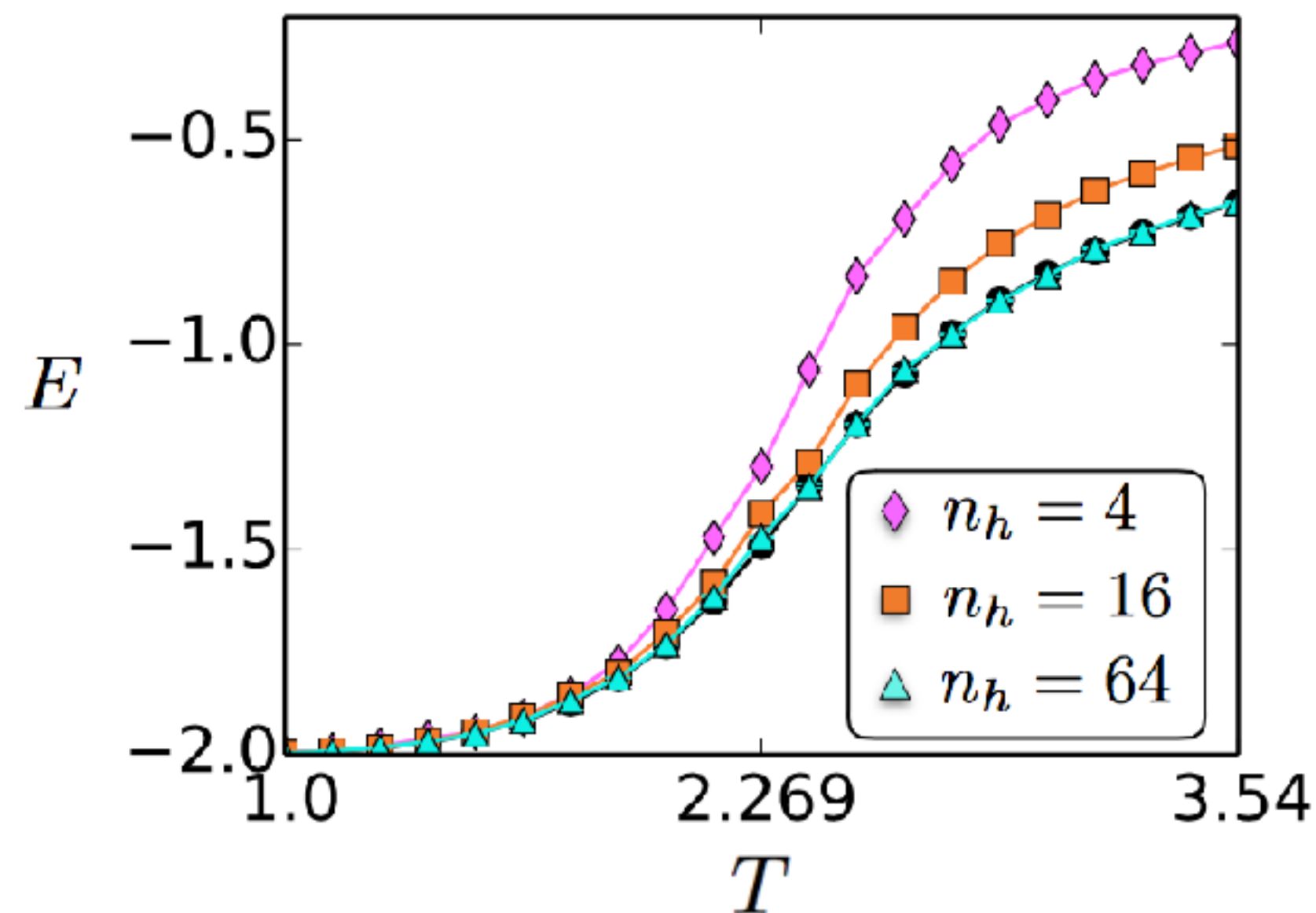
# Learning Thermodynamics of the Ising model

Torlai and RGM, Phys. Rev. B 94, 165134 (2016)

Results from the generative model, after training:    $\langle \mathcal{O} \rangle = \dfrac{1}{N_{\mathrm{MCS}}} \sum_{\mathbf{x}} \mathcal{O}_{\mathbf{x}}$    $\mathbf{x}$ from standard MCMC
= "exact"

$N = 64$



This shows us in this example that the number of hidden units required for accurate generative modelling is approximately the same as the number of hidden units.

See: Chen, Cheng, Xie, Wang, Xiang, Phys. Rev. B 97, 085104 (2018)

Note: The number of measurements required for training & sampling also affects efficiency
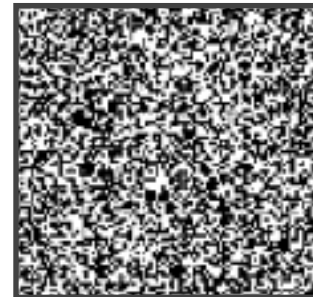
# Learning wavefunctions 1

In the case where the wavefunction is real and positive in a certain basis

$$\psi_\lambda(\mathbf{x}) \propto \sqrt{p_\lambda(\mathbf{x})}$$

Train with samples in the $S^z$ basis



and afterwords calculate estimators from samples produced on the trained machine
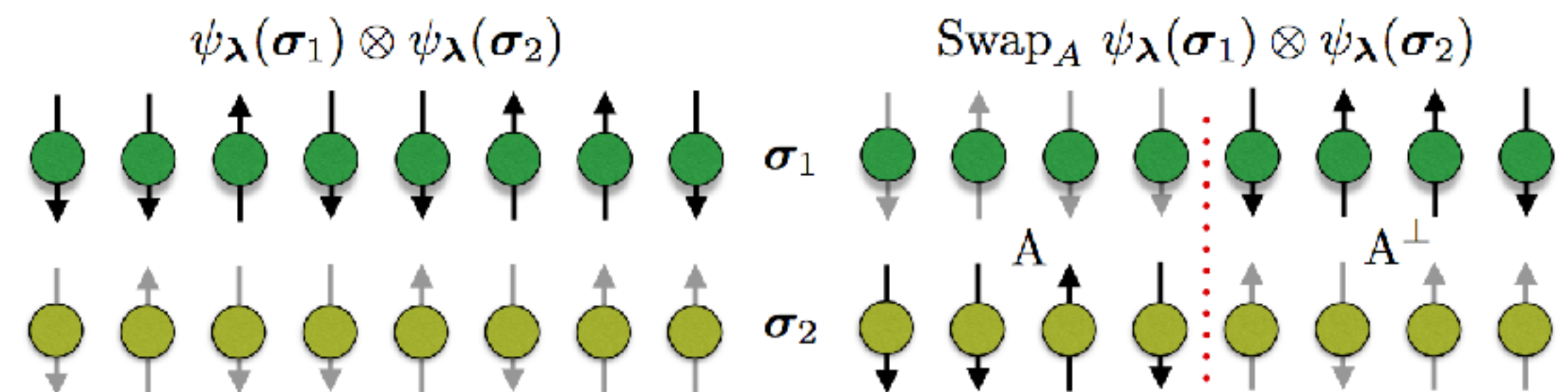
$$\langle \mathcal{O}^{\mathrm{D}} \rangle = \sum_{\mathbf{x}} p_\lambda(\mathbf{x}) \mathcal{O}_{\mathbf{x}}$$

"local" estimator

$$\langle \mathcal{O}^{\mathrm{OD}} \rangle = \sum_{\mathbf{xx'}} \sqrt{p_\lambda(\mathbf{x})}\sqrt{p_\lambda(\mathbf{x'})}\mathcal{O}_{\mathbf{xx'}} = \sum_{\mathbf{x}} p_\lambda(\mathbf{x}) \sum_{\mathbf{x'}} \frac{\sqrt{p_\lambda(\mathbf{x'})}}{\sqrt{p_\lambda(\mathbf{x})}}\mathcal{O}_{\mathbf{xx'}}$$

Entanglement entropy:
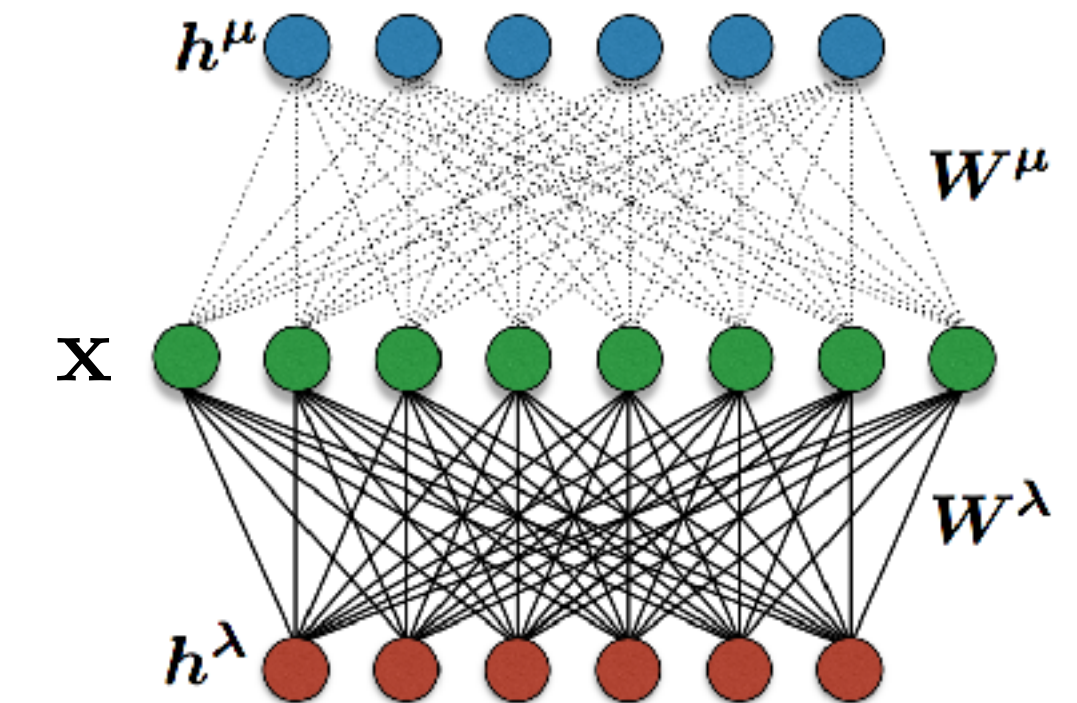
$$S_2(\rho_A) = -\log\left[\mathrm{Tr}(\rho_A^2)\right]$$



$$\psi_{\boldsymbol{\lambda}}(\boldsymbol{\sigma}_1) \otimes \psi_{\boldsymbol{\lambda}}(\boldsymbol{\sigma}_2) \qquad \mathrm{Swap}_A\ \psi_{\boldsymbol{\lambda}}(\boldsymbol{\sigma}_1) \otimes \psi_{\boldsymbol{\lambda}}(\boldsymbol{\sigma}_2)$$

$\boldsymbol{\sigma}_1$

$\boldsymbol{\sigma}_2$

A

$A^\perp$

# Learning wavefunctions 2

For a more generic wavefunction with amplitude and phase, represent both with hidden units

$$\psi_{\lambda,\mu}(\mathbf{x}) \propto \sqrt{p_\lambda(\mathbf{x})}e^{i\phi_\mu(\mathbf{x})}$$



Now, different bases are needed to estimate both the amplitude and phases of the target state.

$$\mathcal{L} = \sum_b^{N_b} \sum_{\mathbf{x}_b} \log|\psi_{\lambda,\mu}(\mathbf{x}_b)|^2$$

$N_b$ = number of bases

$$\{X, X, Z, Z, \dots\}, \{Z, X, X, Z, \dots\}, \{Z, Z, X, X, \dots\},$$

state rotated into basis b with the appropriate unitary

From this, calculate $\nabla_\lambda \mathcal{L}$ and $\nabla_\mu \mathcal{L}$, use stochastic gradient descent, etc.

In practice, training is done in two stages: learning of amplitude first, then optimization of the phase parameters.
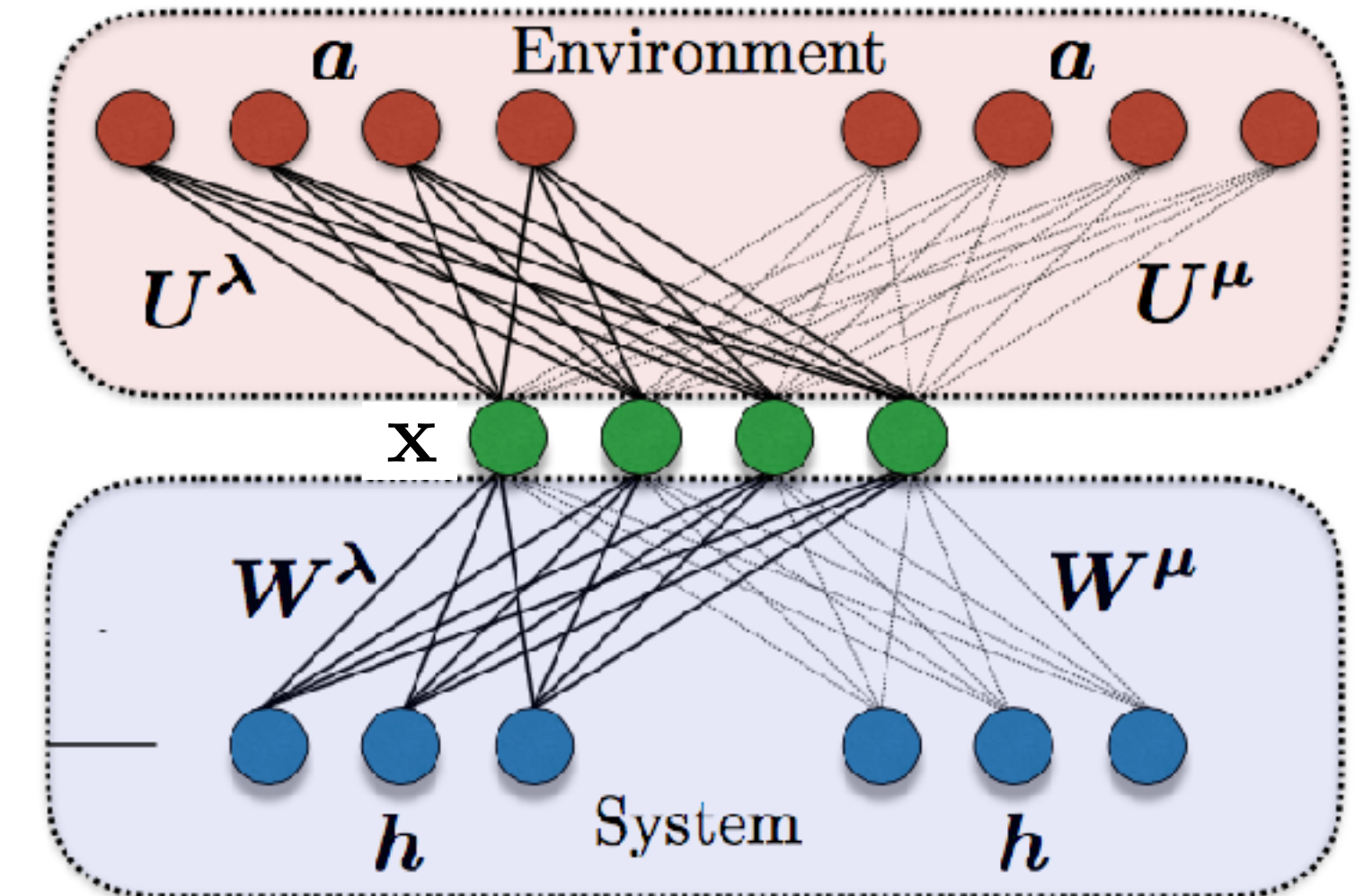
# Learning mixed states

In many experimental setups, quantum states are difficult to isolate, and can be entangled with the environment: one cannot assume *purity*

Can extend our RBM to represent mixed states described by density matrices

$$\psi_{\lambda,\mu}(\mathbf{x}, \mathbf{a}) \propto \sqrt{p_\lambda(\mathbf{x}, \mathbf{a})} e^{i\phi_\mu(\mathbf{x}, \mathbf{a})}$$

$$\rho_{\lambda,\mu}(\mathbf{x}, \mathbf{x}') = \sum_{\mathbf{a}} \psi_{\lambda,\mu}(\mathbf{x}, \mathbf{a}) \psi^*_{\lambda,\mu}(\mathbf{x}', \mathbf{a})$$



Everything we need for *quantum state tomography*

example: Bell state with a global depolarizing channel, 50% error probability

$$\mathcal{F}_{\mathrm{MaxLik}} = 0.9985$$

$$\mathcal{F}_{\mathrm{RBM}} = 0.9992$$

# Rydberg atom arrays



- Neutral atoms (Rb, Sr) are loaded into a lattice formed by an array of optical tweezers

- Atoms can be in their ground state, or an excited state with a large principle quantum number (a Rydberg state).  They form a strongly-interacting system.

- Single-atom resolved fluorescent imaging provides projective measurements

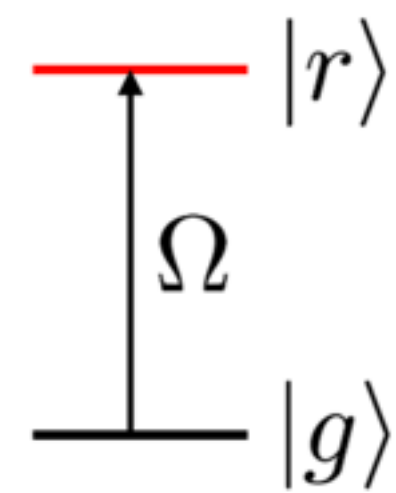- Arrays of atoms are currently used for simulation (groundstates, critical phenomena), solving combinatorial optimization problems

# Rydberg Blockade Hamiltonian

Jaksch, Cirac, Zoller, Rolston, Cote, Lukin, Phys. Rev. Lett. 85, 2208 (2000)
Lukin, Fleischhauer, Cote, Duan, Jaksch, Cirac, Zoller, Phys. Rev. Lett. 87, 037901 (2001)
Fendley, Sengupta, Sachdev, Phys. Rev. B 69, 075106 (2004)

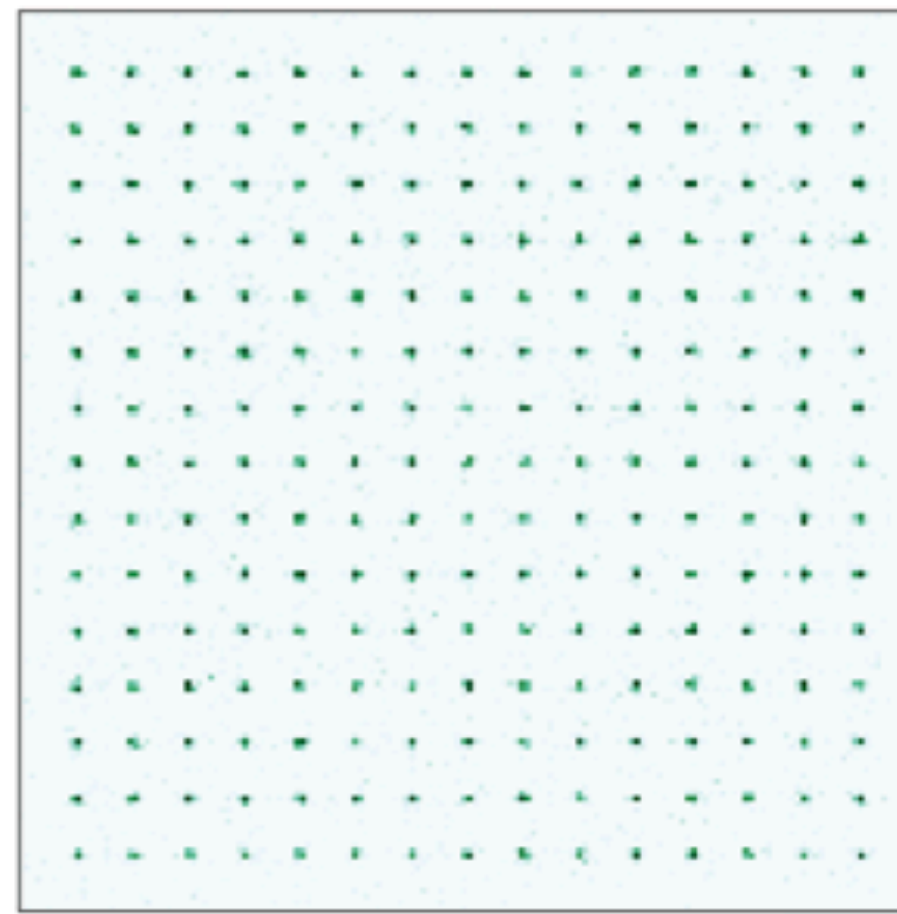$$H = \Omega \sum_i \sigma_i^x - \Delta \sum_i n_i + \sum_{i<j} V_{ij} n_i n_j$$

$$V(R) = \frac{\Omega}{(R/R_b)^6}$$

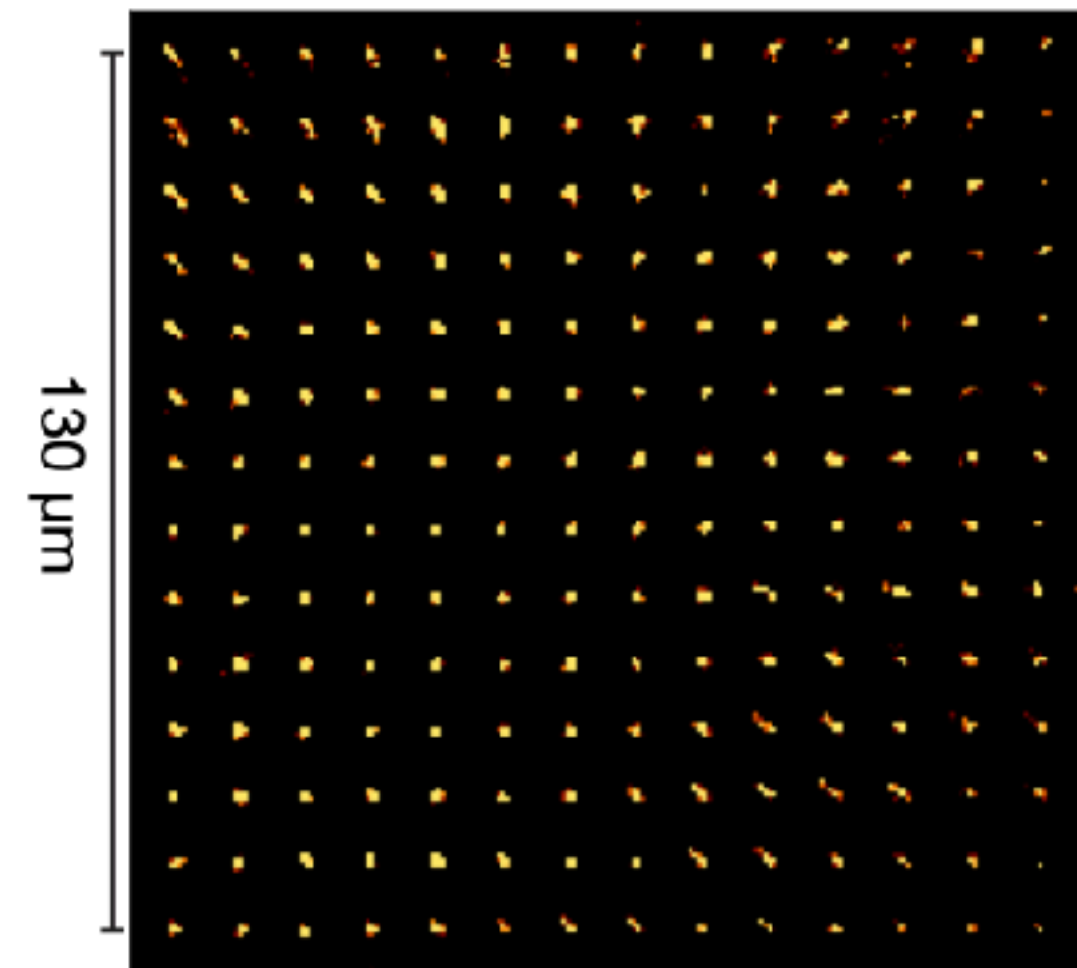$$\sigma^x = |g\rangle\langle r| + |r\rangle\langle g| \qquad n = |r\rangle\langle r|$$

- Two atoms within the blockade radius cannot both be excited into a Rydberg state simultaneously

- Lattice geometry crucially affects physics



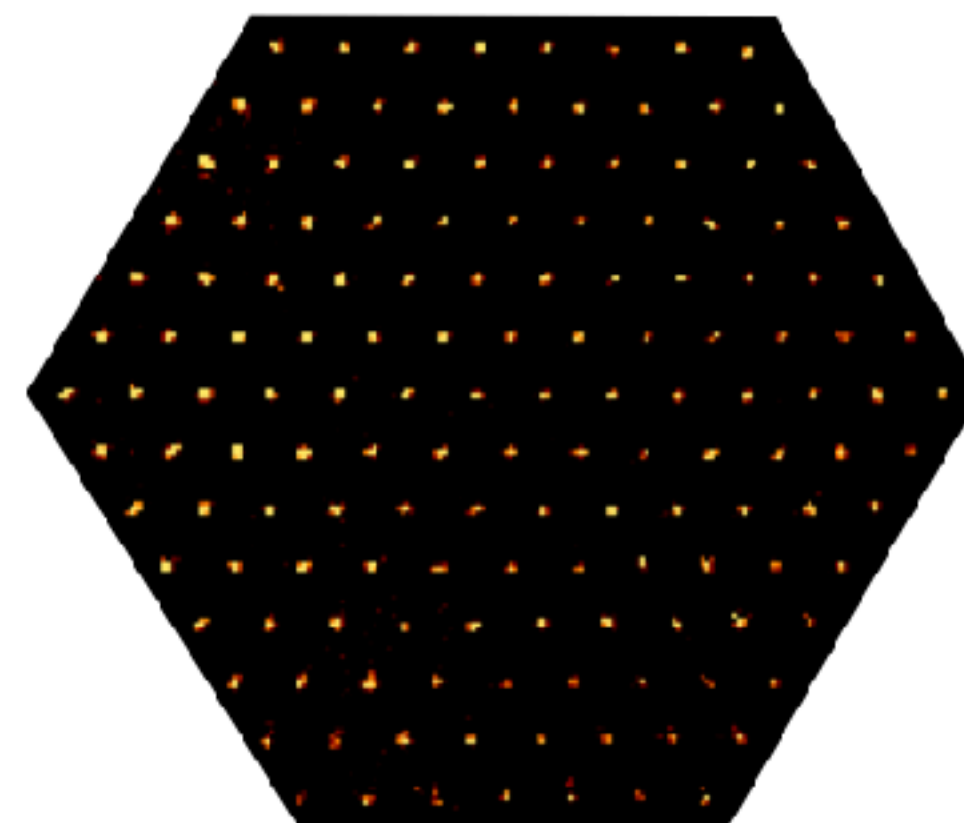Browaeys, Lahaye, Nature Physics 16, 132 (2020)
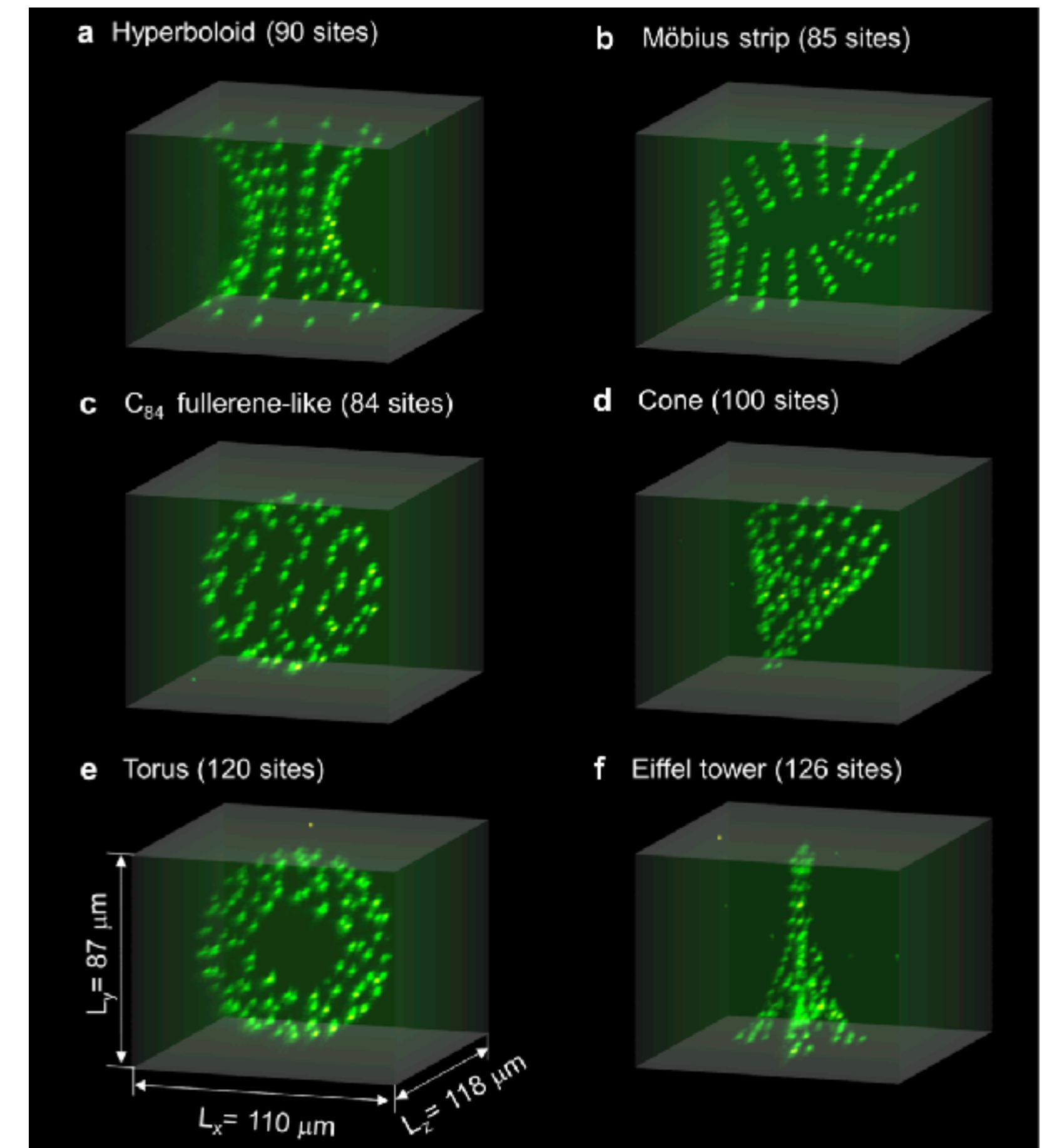
# Experimental lattices



Ebadi et. al. arXiv:2012.12281
Nature 595, 227 (2021)

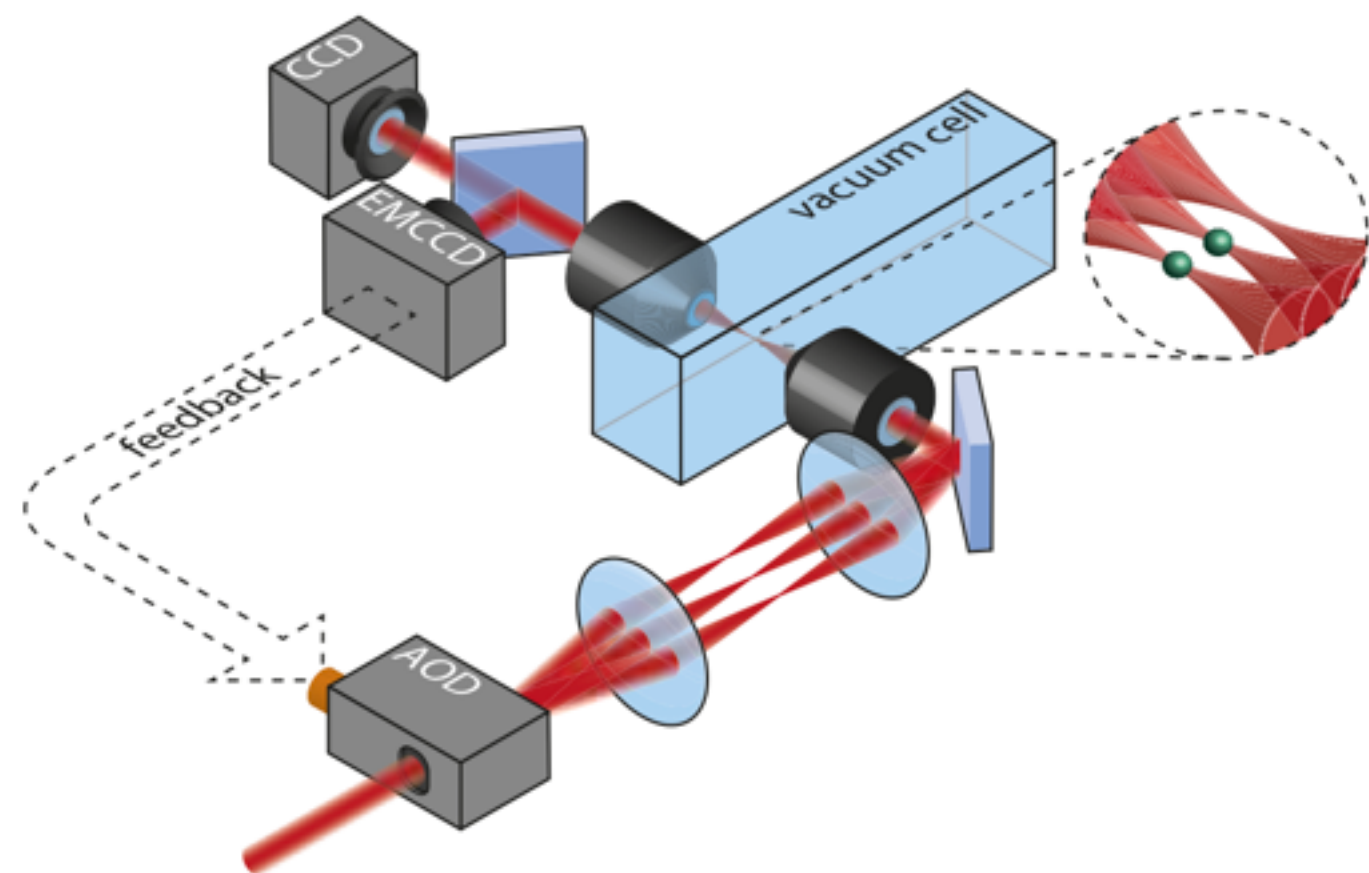Semeghini et. al. arXiv:2104.04119
Science, 374, 1242 (2021)

Scholl et al. arXiv:2012.12268
Nature 595, 233 (2021)

a  Hyperboloid (90 sites)
b  Möbius strip (85 sites)
c  C$_{84}$ fullerene-like (84 sites)
d  Cone (100 sites)
e  Torus (120 sites)
f  Eiffel tower (126 sites)

$L_y = 87$ μm
$L_x = 110$ μm
$L_z = 118$ μm

Barredo, Lienhard, de Léséleuc, Lahaye, Browaeys
Nature 561, (2018)

# Data driven state reconstruction

The availability of high quality projective measurement data allows for state reconstruction, e.g. through the KL divergence or maximum likelihood methods



qubit projective measurement data distributed according to Born rule, $p(\mathbf{x})$
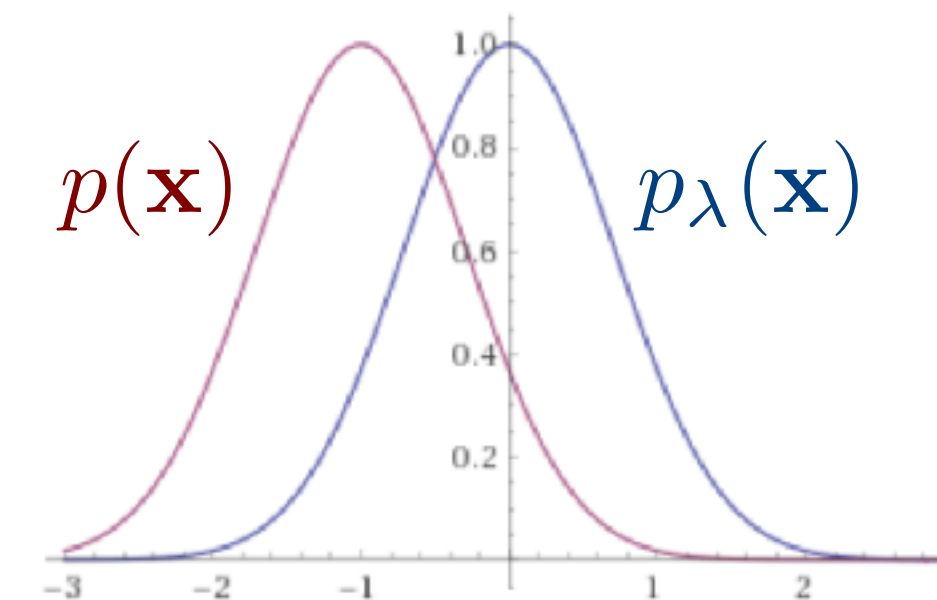
$$\mathbf{x}_1 = (1,0,0,1,1,1,0,0,0,0,\cdots,1)$$
$$\mathbf{x}_2 = (1,1,1,0,1,1,0,1,1,1,\cdots,1)$$
$$\mathbf{x}_3 = (0,1,1,0,0,1,0,1,0,1,\cdots,0)$$

$$\left.\vphantom{\begin{array}{c} \\ \\ \\ \end{array}}\right\} \mathcal{D}$$

Goal: use available data to reconstruct the quantum state using a generative model

$p(\mathbf{x})$        $p_\lambda(\mathbf{x})$

# Rydberg state reconstruction

- one-dimensional chain
- 3000 projective measurements per detuning parameter
- assume positively & purity of the state - use RBM



NN reconstruction

Exact Diagonalization

Stoquastic Hamiltonian:
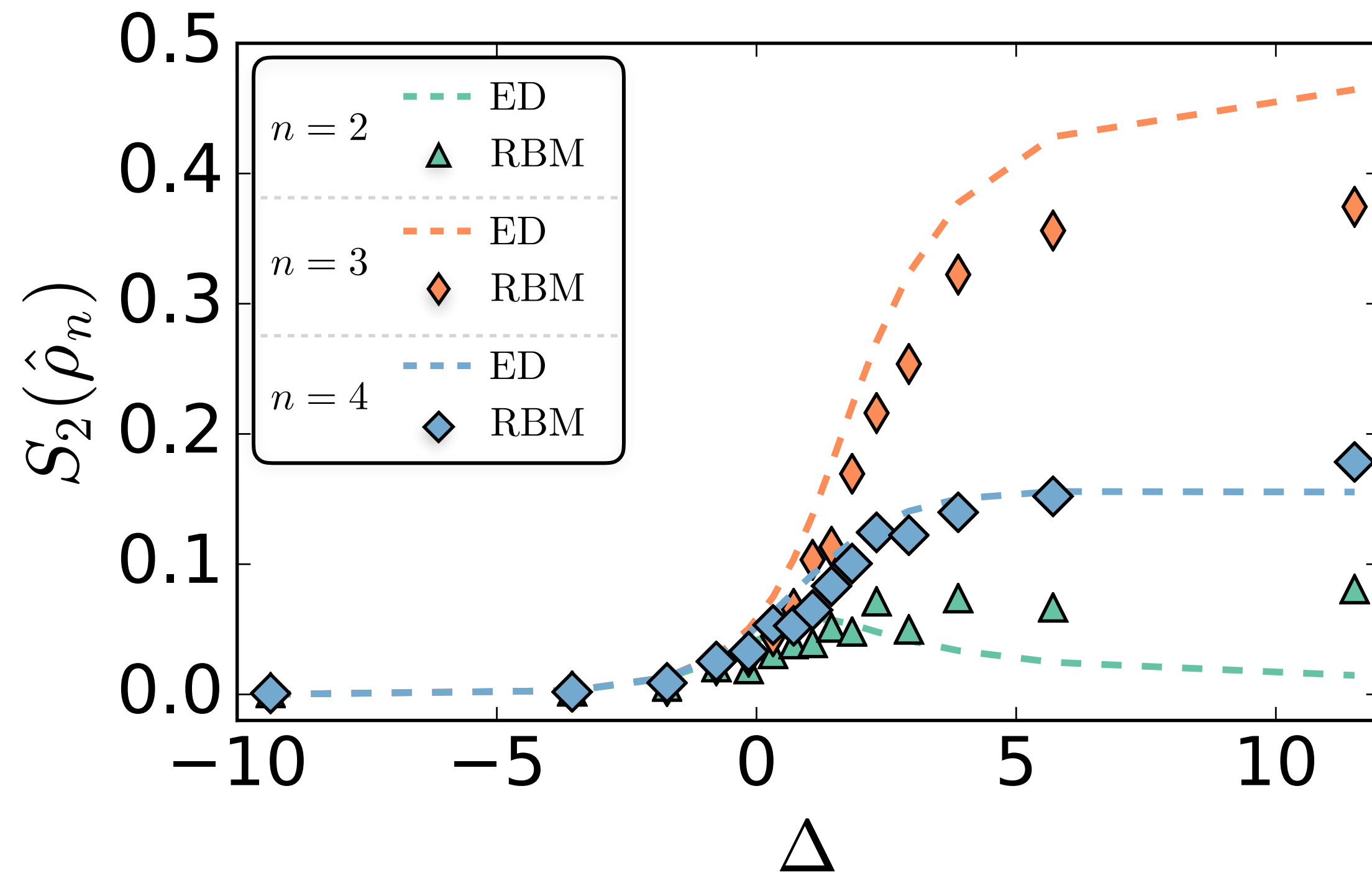
$$\psi_\lambda(z) = \sqrt{p_\lambda(z)}$$

$$\langle \mathcal{O} \rangle = \sum_{zz'} \psi_\lambda(z)\psi_\lambda(z')\mathcal{O}_{zz'}$$

$$= \sum_z \psi_\lambda^2(z) \sum_{z'} \frac{\psi_\lambda(z')}{\psi_\lambda(z)}\mathcal{O}_{zz'}$$

"local" estimator

# Experimental reconstruction



- Second Renyi entropy via SWAP

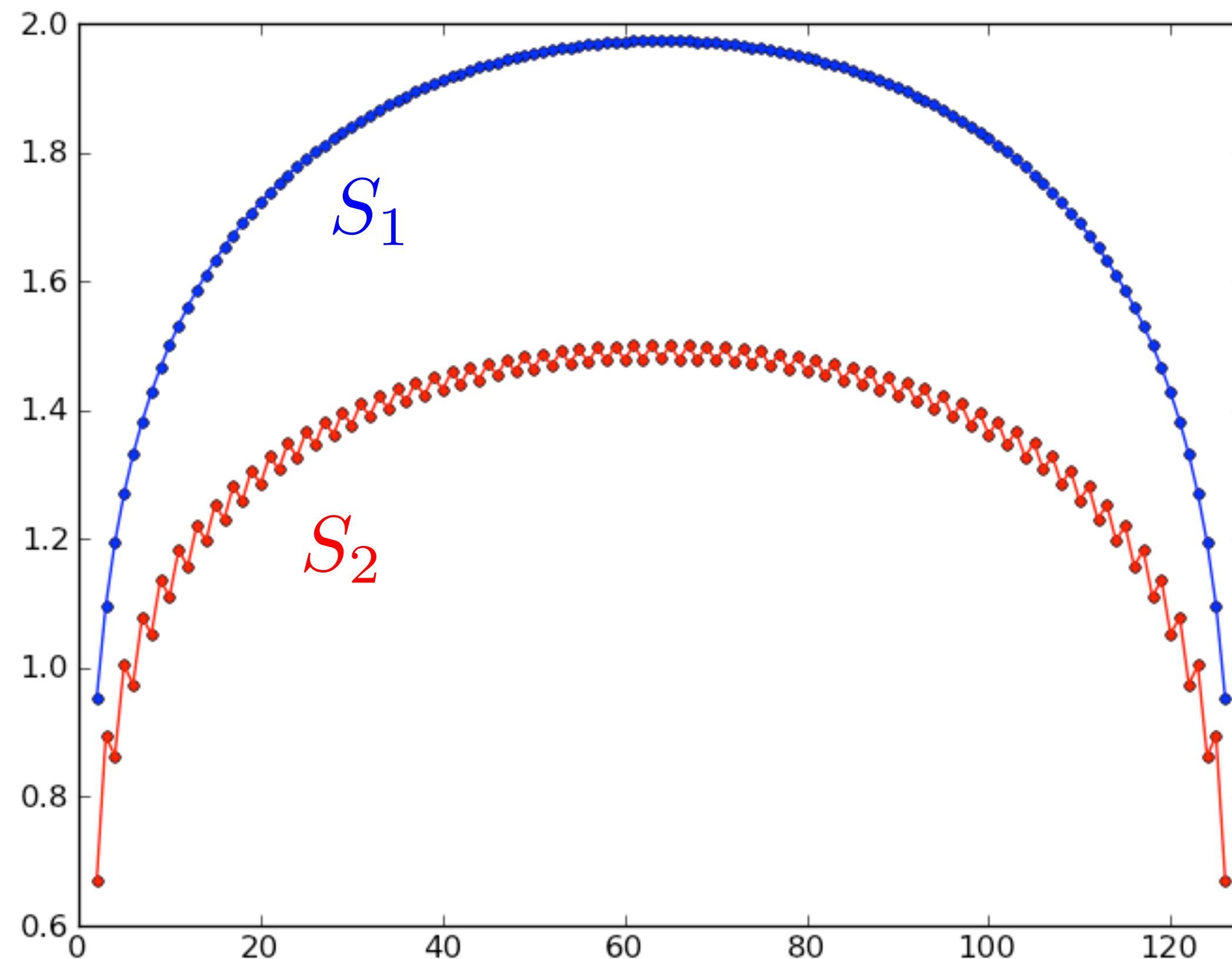$$S_2(\rho_A) = -\log\left[\mathrm{Tr}(\rho_A^2)\right]$$

$$\mathrm{Tr}(\rho_A^2) \Leftrightarrow$$

$$\Leftrightarrow \langle Swap_B \rangle$$

$\langle \sigma^z \rangle$

$\langle \sigma^x \rangle$

ED
Exp
RBM

Legend (left plot):
- ED, RBM $n=2$
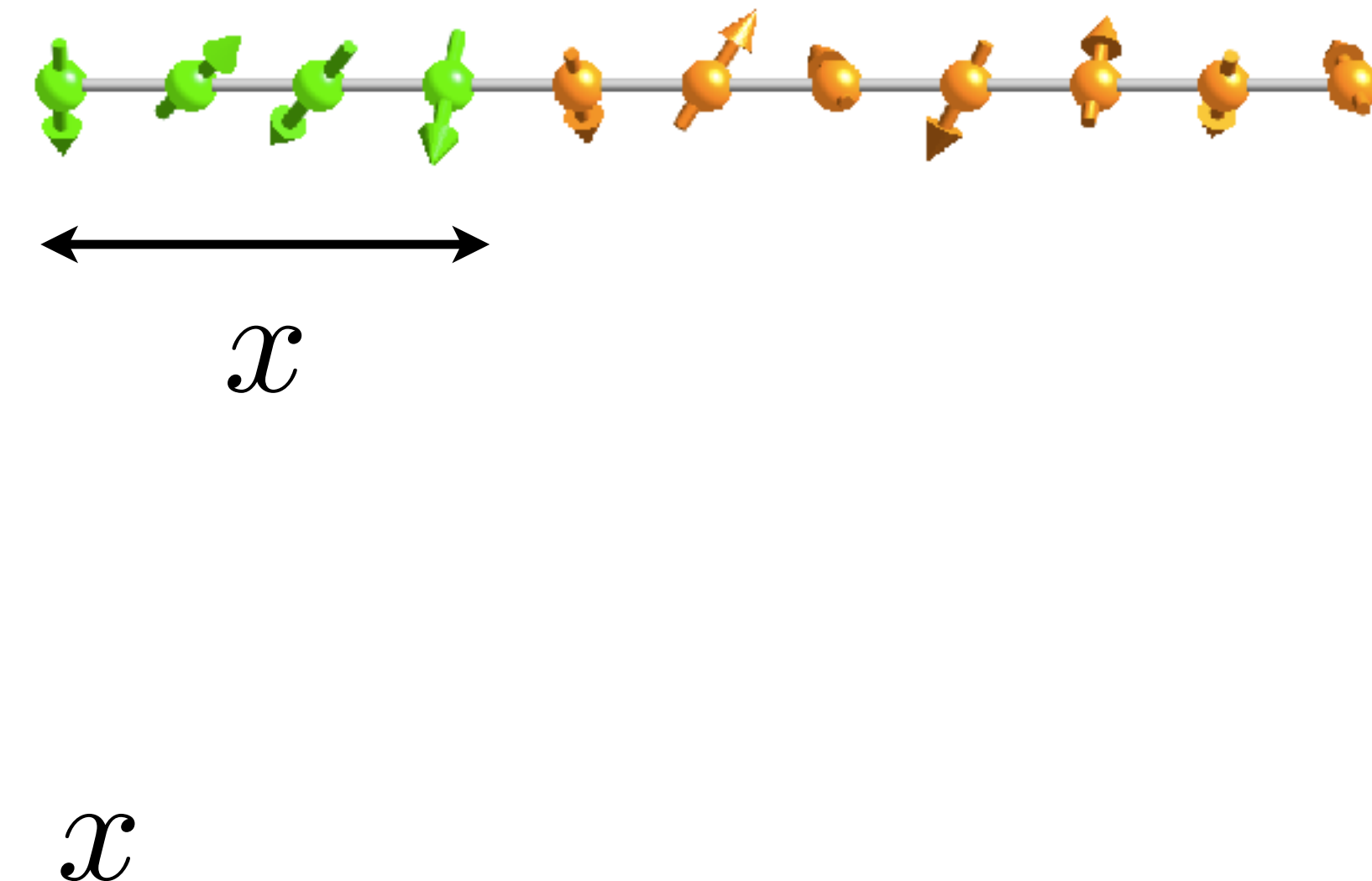- ED, RBM $n=3$
- ED, RBM $n=4$

$S_2(\hat{\rho}_n)$

- **Extracting the central charge**

C. Holzhey, F. Larsen, and F. Wilczek, Nucl. Phys. B424, 443 (1994)
G. Vidal, J. I. Latorre, E. Rico, and A. Kitaev, Phys. Rev. Lett. 90, 227902 (2003)
Calabrese and Cardy, J. Stat. Mech: Theory Exp.  P06002  (2004)

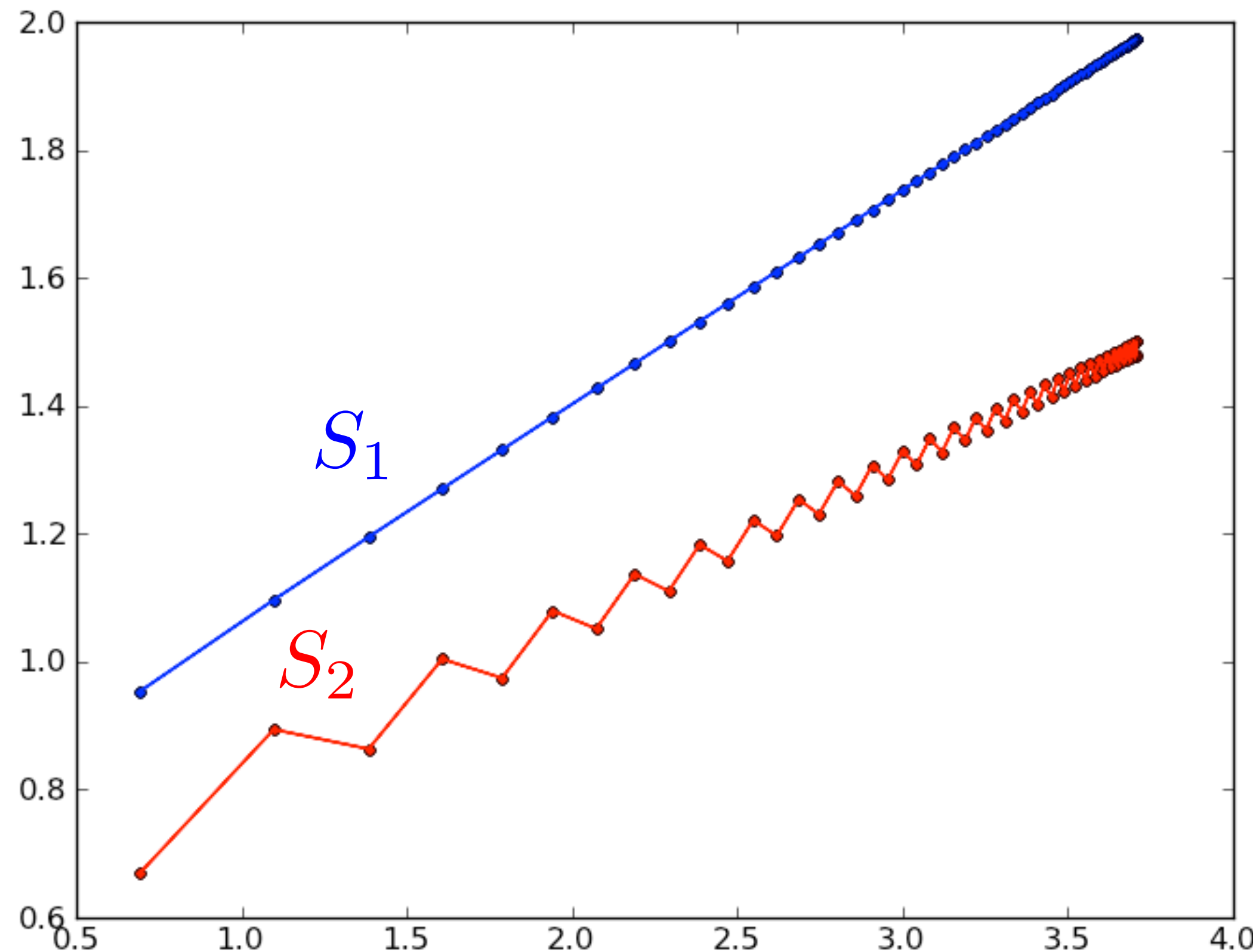$$S_n = \frac{c}{3}\left(1 + \frac{1}{n}\right)\log\left[\frac{L}{\pi a}\sin\frac{\pi x}{L}\right] + \cdots$$



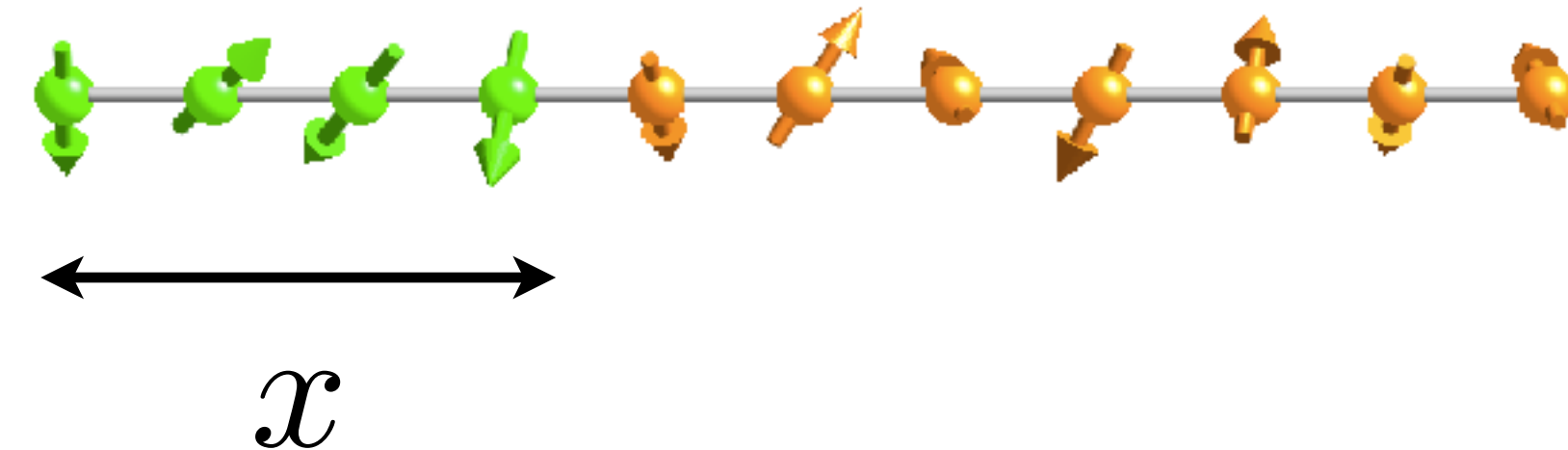$$H = J\sum_{\langle ij\rangle}\mathbf{S}_i\cdot\mathbf{S}_j$$

- Extracting the central charge

C. Holzhey, F. Larsen, and F. Wilczek, Nucl. Phys. B424, 443 (1994)
G. Vidal, J. I. Latorre, E. Rico, and A. Kitaev, Phys. Rev. Lett. 90, 227902 (2003)
Calabrese and Cardy, J. Stat. Mech: Theory Exp. P06002 (2004)

$$S_n = \frac{c}{3}\left(1 + \frac{1}{n}\right)\log\left[\frac{L}{\pi a}\sin\frac{\pi x}{L}\right] + \cdots$$

$$H = J\sum_{\langle ij \rangle}\mathbf{S}_i \cdot \mathbf{S}_j$$

$S_1$

$S_2$

$x$

Calabrese, Campostrini, Essler, Nienhuis
PRL 104, 095701 (2010)

$$\log\left[\frac{L}{\pi}\sin\left(\frac{\pi x}{L}\right)\right]$$