# Neural Network Quantum States and Transfer Learning



Singapore University of Technology and Design



Majulab





# science at a distance

Colloquium

# Machine-Learning for Many-Body Quantum Physics

19 November Thursday, 4.30pm (GMT +8) Giuseppe Carleo EPFL

https://www.quantumlah.org/page/key/upcomingevents

# Plan of the presentation

- 1) Neural Network Quantum States
- 2) Restricted Boltzmann machines
- 3) Restricted Boltzmann machines and transfer learning
- 4) Conclusions and Outlook

Plan of the presentation

- 1) Neural Network Quantum States
- 2) Restricted Boltzmann machines
- 3) Restricted Boltzmann machines and transfer learning
- 4) Conclusions and Outlook

Neural network (NN) can represent probability distributions

In reality they can represent functions ... so why not using neural networks to represent wavefunctions?

Carleo and Troyer (Science 2017) showed that neural networks can be used as a variational ansatz for many-body quantum systems.

This worked open the field of neural-network quantum states NNQS.

- Why now?
- How are NN used in physics?
- How effective NNQS are?
- What is the main working principle for NNQS?
- What architectures for NNQS?

The always more ubiquitous use of NN could not be missed by physicist. Their effectiveness in daily applications is resounding.

Natural language processing, image recognition, games ...



Great progress comes from

- Use of deep networks
- Advances in hardware (GPUs)

Guide experiments

Particle physics and cosmology

Recognize phases of matters

Quantum states and processes tomography

Quantum error correction

Finding the ground state of a Hamiltonian

Finding the steady state of an open quantum system

Perform the time-evolution of a quantum system

Accelerate routines

Zdeborová, L., and F. Krzakala, Adv. Phys. (2016) G Carleo, et al., Reviews of Modern Physics (2019) Universal Approximation Theorem -> NN can represent any well-behaved function. But does it work for the quantum many-body problems we aim to solve?

States described my matrix tensor networks can also be represented by restricted Boltzmann machines (RBM). Chen et al. PRB 2018 This implies that states that can be well represented with tensor networks (see Guo's talk), can also be represented with RBMs.

Restricted Boltzmann machines can represent volume law states. Deng et al. PRX 2017 This implies that there are states which cannot be represented efficiently with tensor networks, but which can be represented efficiently with RBMs. Effectiveness varies significantly on the application and how mature it is.

In particle physics they have been applied for a very long time and they are very effective.

As for neural network quantum states (NNQS), the field is still young and growing

Effectiveness varies significantly on the application and how mature it is.

In particle physics they have been applied for a very long time and they are very effective.

As for neural network quantum states (NNQS), the field is still young and growing

For groundstate, the best comparison is the study of a J1-J2 frustrated system.

Choo et al., Phys. Rev. B 100, 125124 (2019) 10x10 sites



Effectiveness varies significantly on the application and how mature it is.

In particle physics they have been applied for a very long time and they are very effective.

As for neural network quantum states (NNQS), the field is still young and growing

Time evolution in many-body quantum systems.

Carleo and Troyer, Science (2017)



Solid lines results from NNQS

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

We group all the parameters of the network in the vector  $\boldsymbol{\theta}$ 

 $E(\boldsymbol{\theta}) = \frac{\sum_{\boldsymbol{x}, \boldsymbol{x'}} \psi^*(\boldsymbol{x}; \boldsymbol{\theta}) \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \psi(\boldsymbol{x'}; \boldsymbol{\theta})}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''}; \boldsymbol{\theta})|^2}$ 

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

$$E(\boldsymbol{\theta}) = \frac{\sum_{\boldsymbol{x},\boldsymbol{x'}} \psi^*(\boldsymbol{x};\boldsymbol{\theta}) \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \psi(\boldsymbol{x'};\boldsymbol{\theta})}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2}$$
$$= \sum_{\boldsymbol{x}} \left( \frac{|\psi(\boldsymbol{x};\boldsymbol{\theta})|^2}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2} \sum_{\boldsymbol{x'}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \frac{\psi(\boldsymbol{x'};\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})} \right)$$

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

$$E(\boldsymbol{\theta}) = \frac{\sum_{\boldsymbol{x},\boldsymbol{x'}} \psi^*(\boldsymbol{x};\boldsymbol{\theta}) \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \psi(\boldsymbol{x'};\boldsymbol{\theta})}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2}$$
$$= \sum_{\boldsymbol{x}} \left( \frac{|\psi(\boldsymbol{x};\boldsymbol{\theta})|^2}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2} \sum_{\boldsymbol{x'}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \frac{\psi(\boldsymbol{x'};\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})} \right)$$
$$= \sum_{\boldsymbol{x}} p_{\psi}(\boldsymbol{x};\boldsymbol{\theta}) E_{\text{loc}}(\boldsymbol{x};\boldsymbol{\theta})$$

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

$$E(\boldsymbol{\theta}) = \frac{\sum_{\boldsymbol{x},\boldsymbol{x'}} \psi^*(\boldsymbol{x};\boldsymbol{\theta}) \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \psi(\boldsymbol{x'};\boldsymbol{\theta})}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2}$$

$$= \sum_{\boldsymbol{x}} \left( \frac{|\psi(\boldsymbol{x};\boldsymbol{\theta})|^2}{\sum_{\boldsymbol{x''}} |\psi(\boldsymbol{x''};\boldsymbol{\theta})|^2} \sum_{\boldsymbol{x'}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \frac{\psi(\boldsymbol{x'};\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})} \right)$$

$$= \sum_{\boldsymbol{x}} p_{\psi}(\boldsymbol{x};\boldsymbol{\theta}) E_{\text{loc}}(\boldsymbol{x};\boldsymbol{\theta})$$

$$E_{\text{loc}}(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{\boldsymbol{x'}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \frac{\psi(\boldsymbol{x'};\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})}$$

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

$$E(\boldsymbol{\theta}) = \frac{\sum_{\boldsymbol{x},\boldsymbol{x}'} \psi^*(\boldsymbol{x};\boldsymbol{\theta}) \langle \boldsymbol{x} \mid \boldsymbol{H} \mid \boldsymbol{x}' \rangle \psi(\boldsymbol{x}';\boldsymbol{\theta})}{\sum_{\boldsymbol{x}''} |\psi(\boldsymbol{x}'';\boldsymbol{\theta})|^2} \sum_{\boldsymbol{x}'} \langle \boldsymbol{x} \mid \boldsymbol{H} \mid \boldsymbol{x}' \rangle \frac{\psi(\boldsymbol{x}';\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})} \right)$$
From the neural network.
$$= \sum_{\boldsymbol{x}} p_{\psi}(\boldsymbol{x};\boldsymbol{\theta}) E_{\text{loc}}(\boldsymbol{x};\boldsymbol{\theta})$$

$$E_{\text{loc}}(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{\boldsymbol{x}'} \langle \boldsymbol{x} \mid \boldsymbol{H} \mid \boldsymbol{x}' \rangle \frac{\psi(\boldsymbol{x}';\boldsymbol{\theta})}{\psi(\boldsymbol{x};\boldsymbol{\theta})}$$

$$\psi(\boldsymbol{x};\boldsymbol{\theta})$$

Variational principle!

$$E[\psi] = \frac{\langle \psi | H | \psi \rangle}{\langle \psi | \psi \rangle} \ge E_0$$

We group all the parameters of the network in the vector  $\boldsymbol{\theta}$ 



In essence, a variational Monte Carlo method

#### Neural network quantum states: what architectures for NNQS?



Hibat-Allah et al., PRR 2020

#### Neural network quantum states: what architectures for NNQS?



e.g. Metropolis-Hastings, Gibbs

# Plan of the presentation

- 1) Neural Network Quantum States
- 2) Restricted Boltzmann machines
- 3) Restricted Boltzmann machines and transfer learning
- 4) Conclusions and Outlook

Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



Neural networks can be used as a variational ansatz for many-body quantum systems (Carleo and Troyer Science 2017).



$$RBM(\boldsymbol{x}, \boldsymbol{h}; \boldsymbol{W}, \boldsymbol{a}, \boldsymbol{b}) = \frac{1}{Z} e^{\sum_{j} a_{j} x_{j} + \sum_{i} b_{i} h_{i} + \sum_{i,j} x_{j} W_{ji} h_{i}}$$

$$\psi$$

$$p_{\psi} = |\psi|^{2}$$

$$RBM(\boldsymbol{x}; \boldsymbol{W}, \boldsymbol{a}, \boldsymbol{b}) = \frac{1}{Z} \exp\left(\sum_{j} a_{j} x_{j}\right) \prod_{i} 2 \cosh\left(\sum_{j} x_{j} W_{ji} + b_{i}\right)$$
When you know that is real and positive

$$E(\boldsymbol{\theta}) = \sum_{\boldsymbol{x}} p_{\psi}(\boldsymbol{x}; \boldsymbol{\theta}) E_{\text{loc}}(\boldsymbol{x}; \boldsymbol{\theta}) \qquad E_{\text{loc}}(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{\boldsymbol{x'}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x'} \rangle \frac{\psi(\boldsymbol{x'}; \boldsymbol{\theta})}{\psi(\boldsymbol{x}; \boldsymbol{\theta})}.$$

The content of these two boxes allows you do find an approximation for the ground state

# Plan of the presentation

- 1) Neural Network Quantum States
- 2) Restricted Boltzmann machines
- 3) Restricted Boltzmann machines and transfer learning
- 4) Conclusions and Outlook

Singaporean-French collaboration involving SUTD – NUS – Majulab – UCA with 60% physicists and 40% computer scientists.

Remmy Zen, Long My, Ryan Tan, Frederic Hebert, Mario Gattobigio, Christian Miniatura, DP, Stephane Bressan

Singapore

TECHNOLOGY AND DESIGN

France







Zen et al., Physical Review E 2020

Typically transfer learning is used in scenarios like this:

you have a large number of data to train a neural network (for example many photos of cats), and a small amount of data for a different but similar problem (e.g. photos of dogs). Then you can first train the network first with the data and refine it with the cats.



Sample of cats & dogs i

Typically transfer learning is used in scenarios like this:

you have a large number of data to train a neural network (for example many photos of cats), and a small amount of data for a different but similar problem (e.g. photos of dogs). Then you can first train the network first with the data and refine it with the cats.



Sample of cats & dogs images from Kaggle Dataset

Typically transfer learning is used in scenarios like this:

you have a large number of data to train a neural network (for example many photos of cats), and a small amount of data for a different but similar problem (e.g. photos of dogs). Then you can first train the network first with the data and refine it with the cats.



Sample of cats & dogs images from Kaggle Dataset

We can extend this to problem in quantum physics and also **beyond supervised learning**.

Here we aim to work similarly:

- Train a network for a system of a certain size
- Transfer the weights of the network to another network for the description of a larger system.













How do we implement the algorithm?

- To evaluate the energy and the gradients we take 10<sup>4</sup> samples
- To optimize we use a gradient descent algorithm (adaptive learning rate strategies RMSProp and the learning rate of the RMSProp to 0.001
- For the Ising model, the samples are obtained by a single iteration Gibbs sampling.
- For the Heisenberg XXZ model, we use a Metropolis-Hastings sampling with exchange. Strategy to conserve magnetization.
- We use NVIDIA Titan V GPUs

How do we implement the algorithm?

- To evaluate the energy and the gradients we take 10<sup>4</sup> samples
- To optimize we use a gradient descent algorithm (adaptive learning rate strategies RMSProp and the learning rate of the RMSProp to 0.001
- For the Ising model, the samples are obtained by a single iteration Gibbs sampling.
- For the Heisenberg XXZ model, we use a Metropolis-Hastings sampling with exchange. Strategy to conserve magnetization.
- We use NVIDIA Titan V GPUs
- We use twice as many hidden layers as visible ones.

How do we implement the algorithm?

- To evaluate the energy and the gradients we take 10<sup>4</sup> samples
- To optimize we use a gradient descent algorithm (adaptive learning rate strategies RMSProp and the learning rate of the RMSProp to 0.001
- For the Ising model, the samples are obtained by a single iteration Gibbs sampling.
- For the Heisenberg XXZ model, we use a Metropolis-Hastings sampling with exchange. Strategy to conserve magnetization.
- We use NVIDIA Titan V GPUs
- We use twice as many hidden layers as visible ones.
- Stopping criterion. We compute the variance of "local energy"  $E_{\text{loc}}(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{\boldsymbol{x}} \langle \boldsymbol{x} \mid H \mid \boldsymbol{x}' \rangle \frac{\psi(\boldsymbol{x}'; \boldsymbol{\theta})}{\psi(\boldsymbol{x}; \boldsymbol{\theta})}$

$$\sigma_{E_{loc}} = \sqrt{\sum_{\boldsymbol{x}} (E_{loc}(\boldsymbol{x}, \boldsymbol{\theta}) - E(\boldsymbol{\theta}))^2} / \sqrt{N}$$

and we stop when  $\sigma_{E_{loc}}/E(\theta)$  is less than  $\varepsilon_{\sigma} = 0.005$  or when the number of epochs reaches  $\varepsilon_{epoch} = 30,000$ .

We consider two models

$$H_I = -J_I \sum_{\langle l,m\rangle} \sigma_l^z \sigma_m^z - h \sum_l \sigma_l^x \qquad \qquad \text{Ising}$$

$$H_{XXZ} = -J_{XXZ} \sum_{\langle l,m \rangle} \left( \sigma_l^x \sigma_m^x + \sigma_l^y \sigma_m^y + \Delta \sigma_l^z \sigma_m^z \right) \quad \begin{array}{l} \text{Heisenberg} \\ \text{(number conserving)} \end{array}$$



We consider two models

$$H_I = -J_I \sum_{\langle l,m\rangle} \sigma_l^z \sigma_m^z - h \sum_l \sigma_l^x \qquad \qquad \text{Ising}$$

$$H_{XXZ} = -J_{XXZ} \sum_{\langle l,m \rangle} \left( \sigma_l^x \sigma_m^x + \sigma_l^y \sigma_m^y + \Delta \sigma_l^z \sigma_m^z \right) \quad \begin{array}{l} \text{Heisenberg} \\ \text{(number conserving)} \end{array}$$

System sizes 
$$\{4, 8, 16, 32, 64, 128\}$$
 1D  $\{2 \times 2, 4 \times 4, 8 \times 8\}$  2D

We consider two models

$$\begin{split} H_{I} &= -J_{I} \sum_{\langle l,m \rangle} \sigma_{l}^{z} \sigma_{m}^{z} - h \sum_{l} \sigma_{l}^{x} \qquad \text{Ising} \\ H_{XXZ} &= -J_{XXZ} \sum_{\langle l,m \rangle} \left( \sigma_{l}^{x} \sigma_{m}^{x} + \sigma_{l}^{y} \sigma_{m}^{y} + \Delta \sigma_{l}^{z} \sigma_{m}^{z} \right) \qquad \text{Heisenberg} \\ \text{(number conserving)} \end{split}$$

#### We evaluate

- Efficiency: how quickly the ground state is reached
- Effectiveness: how well the ground state is represented

We consider two models

$$\begin{split} H_{I} &= -J_{I} \sum_{\langle l,m \rangle} \sigma_{l}^{z} \sigma_{m}^{z} - h \sum_{l} \sigma_{l}^{x} \qquad \text{Ising} \\ H_{XXZ} &= -J_{XXZ} \sum_{\langle l,m \rangle} \left( \sigma_{l}^{x} \sigma_{m}^{x} + \sigma_{l}^{y} \sigma_{m}^{y} + \Delta \sigma_{l}^{z} \sigma_{m}^{z} \right) \qquad \text{Heisenberg} \\ \text{(number conserving)} \end{split}$$

For efficiency we measure time.

For effectiveness we look at a number of properties and we compare to accurate results from Matrix Product States simulations.

Ferromagnetic order

We compute the **time** required to reach the stopping criterion, adding the time required at previous steps.

We compute the time required to reach the stopping criterion, adding the time required at previous steps.



Let us look at the details of why one transfer learning technique is faster than the other.

We consider only Ising and we compare (1,2) vs (L,2) tiling in the ferromagmetic phase.



We compute the **time** required to reach the stopping criterion, adding the time required at previous steps.



We compute the time required to reach the stopping criterion, adding the time required at previous steps.



How good are the states generated?

How close is the energy from the matrix product states value?

We consider a state for 64 spins and we double its size with different transfer learning techniques and then we do **comparisons to Matrix Product States** results. For a fair comparison, **we consider the same number of epochs for all transfer techniques**.



How good are the states generated?

How close is the energy from the matrix product states value?

We consider a state for 64 spins and we double its size with different transfer learning techniques and then we do **comparisons to Matrix Product States** results. For a fair comparison, **we consider the same number of epochs for all transfer techniques**.



How good are the states generated?

How close is the energy from the matrix product states value?

We consider a state for 64 spins and we double its size with different transfer learning techniques and then we do **comparisons to Matrix Product States** results. For a fair comparison, **we consider the same number of epochs for all transfer techniques**.



Not only the energy, but also the correlations can be predicted well when using the appropriate transfer learning

Ising

Heisenberg

How good are the states generated?

How close is the energy from the matrix product states value?

We consider a state for 64 spins and we double its size with different transfer learning techniques and then we do comparisons to Matrix Product States results. For a fair comparison, we consider the same number of epochs for all transfer techniques.



Ising

# Plan of the presentation

- 1) Neural Network Quantum States
- 2) Restricted Boltzmann machines
- 3) Restricted Boltzmann machines and transfer learning
- 4) Conclusions and Outlook

#### **Conclusions and outlook**

Transfer learning seems to provide an effective and efficient way to study larger quantum many-body systems. What about frustrated Hamiltonians?

Zen et al., Physical Review E 2020





Thank you!

Support



#### Conclusions and outlook

Transfer learning seems to provide an effective and efficient way to study larger quantum many-body systems. What about frustrated Hamiltonians?

Zen et al., Physical Review E 2020

Many lines of research at the interface between machine learning and many-body physics. A lot is yet to be done

- What are the limits of efficiency and effectiveness in finding ground states?
- Are NNQS the way to go for 2D and 3D systems?
- What are the limits of time evolution?
- Can we use more physically-inspired architectures and optimizers so as to improve the performance?
- Can physics teach us how to do better neural networks for any sort of problems?
- Can physics help us to interpret neural networks better?





Thank you!

Support

