

# Variational Neural Annealing

Mohamed Hibat-Allah,<sup>1,2,\*</sup> Estelle M. Inack,<sup>3,1</sup> Roeland Wiersema,<sup>1,2</sup> Roger G. Melko,<sup>2,3</sup> and Juan Carrasquilla<sup>1,2</sup>

<sup>1</sup>*Vector Institute, MaRS Centre, Toronto, Ontario, M5G 1M1, Canada*

<sup>2</sup>*Department of Physics and Astronomy, University of Waterloo, Ontario, N2L 3G1, Canada*

<sup>3</sup>*Perimeter Institute for Theoretical Physics, Waterloo, ON N2L 2Y5, Canada*

(Dated: January 26, 2021)

Many important challenges in science and technology can be cast as optimization problems. When viewed in a statistical physics framework, these can be tackled by simulated annealing, where a gradual cooling procedure helps search for groundstate solutions of a target Hamiltonian. While powerful, simulated annealing is known to have prohibitively slow sampling dynamics when the optimization landscape is rough or glassy. Here we show that by generalizing the target distribution with a parameterized model, an analogous annealing framework based on the variational principle can be used to search for groundstate solutions. Modern autoregressive models such as recurrent neural networks provide ideal parameterizations since they can be exactly sampled without slow dynamics even when the model encodes a rough landscape. We implement this procedure in the classical and quantum settings on several prototypical spin glass Hamiltonians, and find that it significantly outperforms traditional simulated annealing in the asymptotic limit, illustrating the potential power of this yet unexplored route to optimization.

## I. INTRODUCTION

A wide array of complex combinatorial optimization problems can be reformulated as finding the lowest energy configuration of an Ising Hamiltonian of the form [1]:

$$H_{\text{target}} = - \sum_{i < j} J_{ij} \sigma_i \sigma_j - \sum_{i=1}^N h_i \sigma_i, \quad (1)$$

where  $\sigma_i = \pm 1$  are spin variables defined on the  $N$  nodes of a graph. The topology of the graph together with the couplings  $J_{ij}$  and fields  $h_i$  uniquely encode the optimization problem, and its solutions correspond to spin configurations  $\{\sigma_i\}$  that minimize  $H_{\text{target}}$ . While the lowest energy states of certain families of Ising Hamiltonians can be found with modest computational resources, most of these problems are hard to solve and belong to the non-deterministic polynomial time (NP)-hard complexity class [2].

Various heuristics have been used over the years to find approximate solutions to these NP-hard problems. A notable example is simulated annealing (SA) [3], which mirrors the analogous annealing process in materials science and metallurgy where a crystalline solid is heated and then slowly cooled down to its lowest energy and most structurally stable crystal arrangement. In addition to providing a fundamental connection between the thermodynamic behavior of real physical systems and complex optimization problems, simulated annealing has enabled scientific and technological advances with far-reaching implications in areas as diverse as operations research [4], artificial intelligence [5], biology [6], graph theory [7], power systems [8], quantum control [9], circuit design [10] among many others [5]. The paradigm of

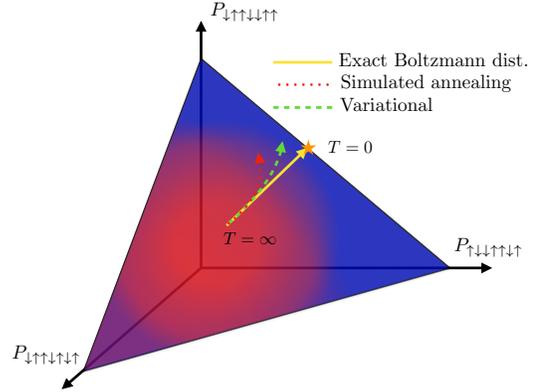


Figure 1. Schematic illustration of the space of probability distributions visited during simulated annealing. An arbitrarily slow SA visits a series of Boltzmann distributions starting at the high temperature (e.g.  $T = \infty$ ) and ending in the  $T = 0$  Boltzmann distribution (continuous yellow line), where a perfect solution to an optimization problem is reached. These solutions are found either at the edge or a corner (for non-degenerate problems) of the standard probabilistic simplex (colored triangle plane). A practical, finite-time SA trajectory (red dotted line), as well as a variational classical annealing trajectory (green dashed line), deviate from the trajectory of exact Boltzmann distributions.

annealing has been so successful that it has inspired intense research into its quantum extension, which requires quantum hardware to anneal the tunneling amplitude, and can be simulated in an analogous way to SA [11, 12].

The SA algorithm explores an optimization problem's energy landscape via a gradual decrease in thermal fluctuations generated by the Metropolis-Hastings algorithm. The procedure stops when all thermal kinetics are removed from the system, at which point the solution to the optimization problem is expected to be found. While an exact solution to the optimization problem is al-

\* mohamed.hibat.allah@uwaterloo.ca

ways attained if the decrease in temperature is arbitrarily slow, a practical implementation of the algorithm must necessarily run on a finite time scale [13]. As a consequence, the annealing algorithm samples a series of effective, quasi-equilibrium distributions close but not exactly equal to the stationary Boltzmann distributions targeted during the annealing [14] (see Fig. 1 for a schematic illustration). This naturally leads to approximate solutions to the optimization problem, whose quality generally depends on the interplay between the problem complexity and the rate at which the temperature is decreased.

In this paper, we offer an alternative route to solving optimization problems of the form of Eq. (1), called *variational neural annealing*. Here, the conventional simulated annealing formulation is substituted with the annealing of a parameterized model. Namely, instead of annealing and approximately sampling the exact Boltzmann distribution, this approach anneals a quasi-equilibrium model, which must be sufficiently expressive and capable of tractable sampling. Fortunately, suitable models have recently been provided by machine learning technology [15–17]. In particular, *neural autoregressive* models combined with variational principles have been shown to accurately describe the equilibrium properties of classical and quantum systems [18–21]. Here, we implement variational neural annealing using autoregressive recurrent neural networks, and show that they offer a powerful alternative to conventional SA and its analogous quantum extension, i.e., simulated quantum annealing (SQA) [11]. This powerful and unexplored route to optimization is schematically illustrated in Fig. 1, where a variational neural annealing trajectory (dashed green arrow) is shown to provide a more accurate approximation to the ideal trajectory (continuous yellow line) than a conventional SA run (dotted red line).

## II. VARIATIONAL CLASSICAL AND QUANTUM ANNEALING

We first consider the variational approach to statistical mechanics [18, 22], where a distribution  $p_{\lambda}(\sigma)$  defined by a set of variational parameters  $\lambda$  is optimized to closely reproduce the equilibrium properties of a system at temperature  $T$ . Following the spirit of SA, we dub our first variational neural annealing algorithm *variational classical annealing* (VCA).

The VCA algorithm searches for the ground state of an optimization problem, encoded in a target Hamiltonian  $H_{\text{target}}$ , by slowly annealing the model’s variational free energy

$$F_{\lambda}(t) = \langle H_{\text{target}} \rangle_{\lambda} - T(t)S_{\text{classical}}(p_{\lambda}), \quad (2)$$

from a high temperature to a low temperature. The quantity  $F_{\lambda}(t)$  provides an upper bound to the true instantaneous free energy and can be used at each annealing stage to update  $\lambda$  through gradient-descent techniques. The brackets  $\langle \dots \rangle_{\lambda}$  denote ensemble averages

taken over the probability  $p_{\lambda}(\sigma)$ . The von Neumann entropy is given by

$$S_{\text{classical}}(p_{\lambda}) = - \sum_{\sigma} p_{\lambda}(\sigma) \log(p_{\lambda}(\sigma)), \quad (3)$$

where the sum runs over all the elements of the state space  $\{\sigma\}$ . In our setting, the temperature is decreased from an initial value  $T_0$  to 0 using a linear schedule function  $T(t) = T_0(1 - t)$ , where  $t \in [0, 1]$ , which follows closely the traditional implementation of SA.

In order for VCA to succeed, we require parameterized models that enable the estimation of entropy, Eq. (3), without incurring expensive calculations of the partition function. In addition, we anticipate that hard optimization problems will induce a complex energy landscape into the parameterized models and an ensuing slowdown of their sampling via Markov chain Monte Carlo. These issues preclude un-normalized models such as restricted Boltzmann machines, where sampling relies on Markov chains and whose partition function is intractable to evaluate [23]. Instead, we implement VCA using recurrent neural networks (RNNs) [20, 21], whose autoregressive nature enables statistical averages over exact samples  $\sigma$  drawn from  $p_{\lambda}(\sigma)$ . Since RNNs are normalized by construction, these samples naturally allow the estimation of the entropy in Eq. (3). We provide a detailed description of the RNN in Methods Sec. V A.

The VCA algorithm, summarized in Fig. 2(a), performs a warm-up step which brings a randomly initialized distribution  $p_{\lambda}(\sigma)$  to an approximate equilibrium state with free energy  $F_{\lambda}(t = 0)$  via  $N_{\text{warmup}}$  gradient descent steps. At each step  $t$ , we reduce the temperature of the system from  $T(t)$  to  $T(t + \delta t)$  and apply  $N_{\text{train}}$  gradient descent steps to re-equilibrate the model. A critical ingredient to the success of VCA is that the variational parameters optimized at temperature  $T(t)$  are reused at temperature  $T(t + \delta t)$  to ensure that the model’s distribution is always near its instantaneous equilibrium state. Repeating the last two steps  $N_{\text{annealing}}$  times, we reach temperature  $T(1) = 0$ , which is the end of the annealing protocol. Here the distribution  $p_{\lambda}(\sigma)$  is expected to assign high probability to configurations  $\sigma$  that solve the optimization problem. Likewise, the residual entropy Eq. (3) at  $T(1) = 0$  provides a heuristic approach to count the number of solutions to the problem Hamiltonian [18]. Further algorithmic details are provided in Methods Sec. V B.

Simulated annealing provides a powerful heuristic for the solution of hard optimization problems by harnessing thermal fluctuations. Inspired by the latter, the advent of commercially available quantum devices [24] has enabled the analogous concept of quantum annealing [25], where the solution to an optimization problem is performed by harnessing quantum fluctuations. In quantum annealing, the search for the ground state of Eq. (1) is performed at  $T = 0$ , by supplementing the target Hamiltonian with a quantum mechanical kinetic (or “driving”) term,

$$\hat{H}(t) = \hat{H}_{\text{target}} + f(t)\hat{H}_D, \quad (4)$$

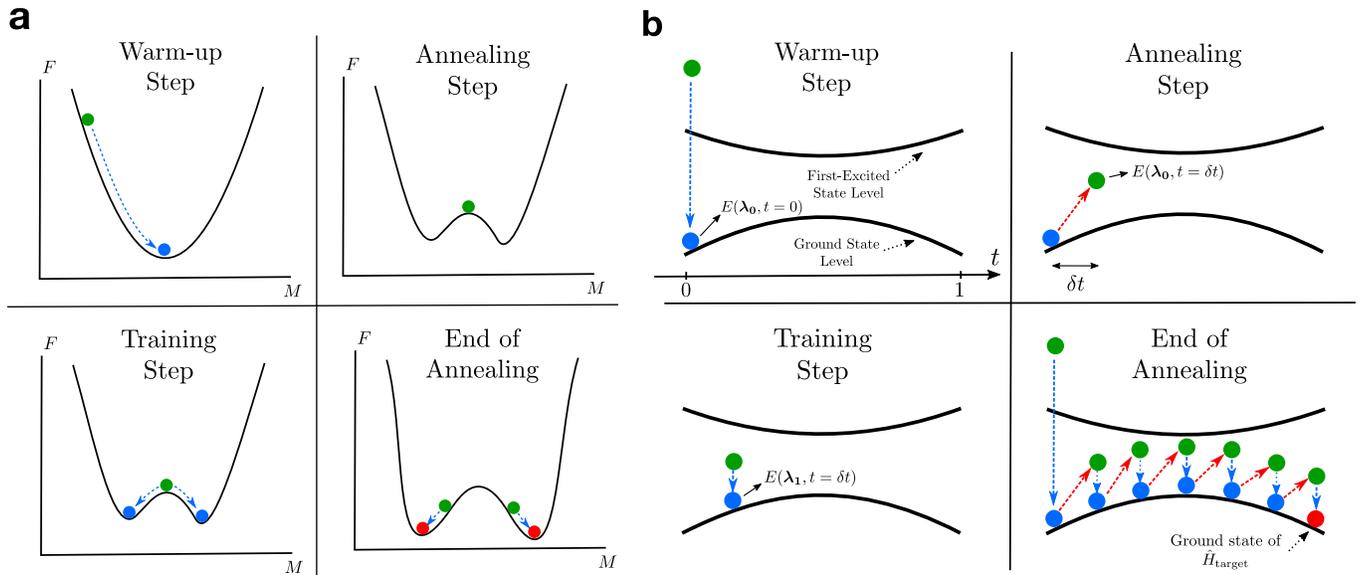


Figure 2. Variational neural annealing protocols. (a) The variational classical annealing (VCA) algorithm steps. A warm-up step brings the initialized variational state (green dot) close to the minimum of the free energy (cyan dot) at a given value of the order parameter  $M$ . This step is followed by an annealing and a training step that brings the variational state back to the new free energy minimum. Repeating the last two steps until  $T(t=1) = 0$  (red dots) produces approximate solutions to  $H_{\text{target}}$  if the protocol is conducted slowly enough. This schematic illustration corresponds to annealing through a continuous phase transition with an order parameter  $M$ . (b) Variational quantum annealing (VQA). VQA includes a warm-up step, followed by an annealing and a training step, which brings the variational energy (green dot) closer to the new a ground state energy (cyan dot). We loop over the previous two steps until reaching the target ground state of  $\hat{H}_{\text{target}}$  (red dot) if annealing is performed slowly enough.

where  $H_{\text{target}}$  in Eq. (1) is promoted to a quantum mechanical Hamiltonian  $\hat{H}_{\text{target}}$ .

Quantum annealing algorithms typically start with a dominant driving term  $\hat{H}_D \gg \hat{H}_{\text{target}}$  chosen so that the ground state of  $\hat{H}(0)$  is easy to prepare. When the strength of the driving term is subsequently reduced (typically adiabatically) using a schedule function  $f(t)$ , the system is annealed to the ground state of  $\hat{H}_{\text{target}}$ . In analogy to its thermal counterpart, SQA emulates this process on classical computers using quantum Monte Carlo methods [11].

Here, we leverage the variational principle of quantum mechanics and devise a strategy that emulates quantum annealing variationally. We dub our second variational neural annealing algorithm *variational quantum annealing* (VQA). The latter is based on the variational Monte Carlo (VMC) algorithm, whose goal is to simulate the equilibrium properties of quantum systems at zero temperature (see Methods Sec. VC). In VMC, the ground state of a Hamiltonian  $\hat{H}$  is modeled through an ansatz  $|\Psi_{\lambda}\rangle$  endowed with parameters  $\lambda$ . The variational principle guarantees that the energy  $\langle \Psi_{\lambda} | \hat{H} | \Psi_{\lambda} \rangle$  is an upper bound to the ground state energy of  $\hat{H}$ , which we use to define a time-dependent objective function  $E(\lambda, t) \equiv \langle \hat{H}(t) \rangle_{\lambda} = \langle \Psi_{\lambda} | \hat{H}(t) | \Psi_{\lambda} \rangle$  to optimize the parameters  $\lambda$ .

The VQA setup, graphically summarized in Fig. 2(b),

applies  $N_{\text{warmup}}$  gradient descent steps to minimize  $E(\lambda, t=0)$ , which brings  $|\Psi_{\lambda}\rangle$  close to the ground state of  $\hat{H}(0)$ . Setting  $t = \delta t$  while keeping the parameters  $\lambda_0$  fixed results in a variational energy  $E(\lambda_0, t = \delta t)$ . A set of  $N_{\text{train}}$  gradient descent steps bring the ansatz closer to the new instantaneous ground state, which results in a variational energy  $E(\lambda_1, t = \delta t)$ . The variational parameters optimized at time step  $t$  are reused at time  $t + \delta t$ , which promotes the computational adiabaticity of the protocol (see Appendix. A). We repeat the annealing and training steps  $N_{\text{annealing}}$  times on a linear schedule ( $f(t) = 1 - t$  with  $t \in [0, 1]$ ) until  $t = 1$ , at which point the system should solve the optimization problem (red dot in Fig. 2(b)). We note that in our simulations, no training steps are taken at  $t = 1$ . Finally, similarly to VCA, we choose normalized RNN wave functions [20, 21] as ansätze, giving the VQA algorithm access to exact Monte Carlo samples.

To gain theoretical insight on the principles behind a successful VQA simulation, we derive a variational version of the adiabatic theorem [26]. Starting from a set of assumptions, such as the convexity of the energy landscape in the warm-up phase and close to convergence during annealing, as well as the absence of noise in the energy gradients, we provide a bound on the total number of gradient descent steps  $N_{\text{steps}}$  that guarantees the adiabaticity of the VQA algorithm as well as a success probability of solving the optimization problem  $P_{\text{success}} > 1 - \epsilon$ .

Here,  $\epsilon$  is an upper bound on the overlap between the variational wave function and the excited states of the Hamiltonian  $\hat{H}(t)$ , i.e.,  $|\langle \Psi_{\perp}(t) | \Psi_{\lambda} \rangle|^2 < \epsilon$ . We show that  $N_{\text{steps}}$  can be bounded as (see Appendix. B):

$$\mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon \min_{\{t_n\}}(g(t_n))}\right) \leq N_{\text{steps}} \leq \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon^2 \min_{\{t_n\}}(g(t_n))^2}\right). \quad (5)$$

The function  $g(t)$  is the energy gap between the first excited state and the ground state of the instantaneous Hamiltonian  $\hat{H}(t)$ ,  $N$  is the system size, and the set of times  $\{t_n\}$  is defined in Appendix. B. As expected for hard optimization problems, the minimum gap typically decreases exponentially with system size  $N$ , which dominates the computational complexity of a VQA simulation, but in cases where the minimum gap scales as the inverse of a polynomial in  $N$ , then the number of steps  $N_{\text{steps}}$  is also polynomial in  $N$ .

### III. RESULTS

#### A. Annealing on random Ising chains

We now proceed to evaluate the power of VCA and VQA. As a first benchmark, we consider the task of solving for the ground state the one-dimensional (1D) Ising Hamiltonian with random couplings  $J_{i,i+1}$ ,

$$H_{\text{target}} = - \sum_{i=1}^{N-1} J_{i,i+1} \sigma_i \sigma_{i+1}. \quad (6)$$

First, we examine  $J_{i,i+1}$  sampled from a uniform distribution in the interval  $[0, 1)$ . Here, the ground state configuration is given either by all spins up or down, and the ground state energy is known exactly, i.e.,  $E_G = - \sum_{i=1}^{N-1} J_{i,i+1}$  [27].

We use a tensorized RNN ansatz without weight sharing for both VCA and VQA (see Methods Sec. VA). We consider system sizes  $N = 32, 64, 128$  and  $N_{\text{train}} = 5$ , which suffices to achieve accurate solutions. For VQA, we use a one-body driving term  $\hat{H}_D = -\Gamma_0 \sum_{i=1}^N \hat{\sigma}_i^x$ , where  $\hat{\sigma}_i^{x,y,z}$  are Pauli matrices acting on site  $i$ . To quantify the performance of the algorithms, we use the residual energy [11],

$$\epsilon_{\text{res}} = [\langle H_{\text{target}} \rangle_{\text{av}} - E_G]_{\text{dis}}, \quad (7)$$

where  $E_G$  is the exact ground state energy of  $H_{\text{target}}$ . We use the arithmetic mean for statistical averages  $\langle \dots \rangle_{\text{av}}$  over samples from the models. For VCA it means that  $\langle H_{\text{target}} \rangle_{\text{av}} \approx \langle H_{\text{target}} \rangle_{\lambda}$ , while for VQA the target Hamiltonian is promoted to  $\hat{H}_{\text{target}} = - \sum_{i=1}^{N-1} J_{i,i+1} \hat{\sigma}_i^z \hat{\sigma}_{i+1}^z$  and  $\langle H_{\text{target}} \rangle_{\text{av}} \approx \langle \hat{H}_{\text{target}} \rangle_{\lambda}$ . We consider the typical (geometric) mean for averaging over instances of the target Hamiltonian, i.e.,  $[\dots]_{\text{dis}} = \exp(\langle \ln(\dots) \rangle_{\text{av}})$ . The average in the argument of the exponential stands for arithmetic mean over different realizations of the couplings.

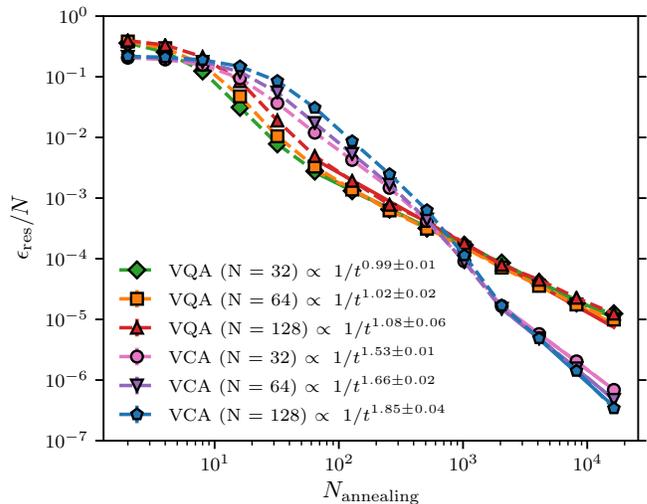


Figure 3. Variational neural annealing on a random Ising chain. Here we represent the residual energy per site  $\epsilon_{\text{res}}/N$  vs the number of annealing steps  $N_{\text{annealing}}$  for both VQA and VCA. The system sizes are  $N = 32, 64, 128$ . We use random positive couplings  $J_{i,i+1} \in [0, 1)$  (see text for more details). The error bars represent the one s.d. statistical uncertainty calculated over different disorder realizations [28].

We take advantage of the autoregressive nature of the RNN and sample  $10^6$  configurations at the end of the annealing, which allows us to accurately estimate the model's arithmetic mean. The typical mean is taken over 25 instances of  $H_{\text{target}}$ .

In Fig. 3 we report the residual energies per site against the number of annealing steps  $N_{\text{annealing}}$ . As expected, the residual energy is a decreasing function of  $N_{\text{annealing}}$ , which underlines the importance of adiabaticity and annealing in our setting. In our examples, we observe that the decrease of the residual energy of VCA and VQA is consistent with a power-law decay for a large number of annealing steps. Whereas VCA's decay exponent is in the interval 1.5 – 1.9, the VQA exponent is about 0.9 – 1.1. These exponents suggest an asymptotic speed-up compared to SA and coherent quantum annealing, where the residual energies follow a logarithmic law [29]. Contrary to the observations in Ref. [29] where quantum annealing was found superior to SA, VCA finds an average residual energy an order of magnitude more accurate than VQA for a large number of annealing steps.

Finally, we note that the exponents provided above are not expected to be universal and are a priori sensitive to the hyperparameters of the algorithms, e.g., learning rate, model choice, number of training steps, optimizer, etc. Appendix. C provides a summary of the hyperparameters used in our work. Additional illustrations of the adiabaticity of VCA and VQA, as well as of the annealing results for a chain with  $J_{i,i+1}$  uniformly sampled from the discrete set  $\{-1, +1\}$ , are provided in Appendix. A.

## B. Edwards-Anderson model

We now consider the two-dimensional (2D) Edwards-Anderson (EA) model, which is a prototypical spin glass arranged on a square lattice with nearest neighbor random interactions. The problem of finding ground states of the model has been studied experimentally [12] and numerically [11] from the annealing perspective, as well as theoretically [2] from the computational complexity perspective. The EA model with open boundary conditions is given by

$$H_{\text{target}} = - \sum_{\langle i,j \rangle} J_{ij} \sigma_i \sigma_j, \quad (8)$$

where  $\langle i,j \rangle$  denote nearest neighbors. The couplings  $J_{ij}$  are drawn from a uniform distribution in the interval  $[-1, 1)$ . In the absence of a longitudinal field, for which solving the EA model is NP-hard, the ground state can be found in polynomial time [2]. To find the exact ground state of each random realization, we use the spin-glass server [30].

We use a 2D tensorized RNN ansatz without weight sharing for the variational protocols (see Methods Sec. V A). For VQA, we use a one-body driving term  $\hat{H}_D = -\Gamma_0 \sum_{i=1}^N \hat{\sigma}_i^x$ . Fig. 4(a) shows the annealing results obtained on a system size  $N = 10 \times 10$  spins. VCA outperforms VQA and in the adiabatic, long-time annealing regime, it produces solutions three orders of magnitude more accurate on average than VQA. In addition, we investigate the performance of VQA supplemented with a fictitious Shannon information entropy [21] term that mimics thermal relaxation effects observed in quantum annealing hardware [31]. This form of regularized VQA, here labelled (RVQA), is described by a pseudo free energy cost function  $\tilde{F}_\lambda(t) = \langle \hat{H}(t) \rangle_\lambda - T(t) S_{\text{classical}}(|\Psi_\lambda|^2)$ . As in VCA, the pseudo entropy term  $S_{\text{classical}}(|\Psi_\lambda|^2)$  at  $f(1) = 0$  provides a heuristic approach to count the number of solutions to  $H_{\text{target}}$  for VQA and RVQA. The results in Fig. 4(a) do show an amelioration of the VQA performance, including changing a saturating dynamics at large  $N_{\text{annealing}}$  to a power-law like behavior. However, it appears to be insufficient to compete with the VCA scaling (see exponents in Fig. 4(a)). This observation suggests the superiority of a thermally driven variational emulation of annealing over a purely quantum one for this example.

To further scrutinize the relevance of the annealing effects in VCA, we also consider VCA with zero thermal fluctuations, i.e., setting  $T_0 = 0$ . Because of its intimate relation to the classical-quantum optimization (CQO) methods of Refs. [32–34], we refer to this setting as CQO. Fig. 4(a) shows that CQO takes about  $10^3$  training steps to reach accuracies nearing 1%. The accuracy does not further improve upon additional training up to  $10^5$  gradient steps, which indicates that CQO is prone to getting stuck in local minima. In comparison, VCA and VQA offer solutions orders of magnitude more ac-

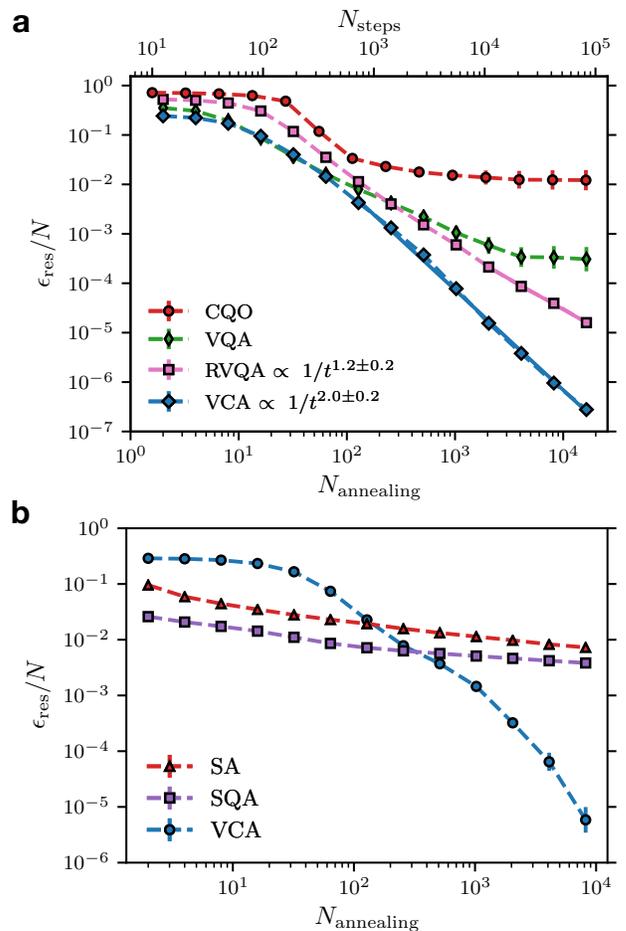


Figure 4. Benchmarking the two-dimensional Edwards-Anderson spin glass. (a) A comparison between VCA, VQA, RVQA, and CQO on a  $10 \times 10$  lattice by plotting the residual energy per site vs  $N_{\text{annealing}}$ . For CQO, we report the residual energy per site vs the number of optimization steps  $N_{\text{steps}}$ . (b) Comparison between SA, SQA with  $P = 20$  trotter slices, and VCA using a 2D tensorized RNN ansatz on a  $40 \times 40$  lattice. The annealing speed is the same for SA, SQA and VCA.

curate on average for a large number of annealing steps, highlighting the importance of annealing in tackling optimization problems.

Since VCA displays the best performance in the previous benchmarks, we use it to demonstrate its capabilities on a  $40 \times 40$  spin system. For comparison, we use SA as well as SQA. The SQA simulation uses the path-integral Monte Carlo method [11] with  $P = 20$  trotter slices, and we report averages over energies across all trotter slices, for each realization of randomness (see Methods Sec. V D). In addition, we average the energy obtained after 25 annealing runs on every instance of randomness for SA and SQA. To average over Hamiltonian instances, we use the typical mean over 25 different realizations for the three annealing methods. The results are shown in Fig. 4(b), where we present the residual

energies per site against the number of annealing steps  $N_{\text{annealing}}$ , which is set so that the speed of annealing is the same for SA, SQA and VCA. We first note that our results confirm the qualitative behavior of SA and SQA in Refs. [11, 35]. While SA and SQA produce lower residual energy solutions than VCA for small  $N_{\text{annealing}}$ , we observe that VCA achieves residual energies about three orders of magnitude smaller than SQA and SA for a large number of annealing steps. Notably, the rate at which the residual energy improves with increasing  $N_{\text{annealing}}$  is significantly higher for VCA compared to SQA and SA even at relatively small number of annealing steps.

### C. Fully-connected spin glasses

We now focus our attention on fully-connected spin glasses [2, 36]. We first focus on the Sherrington-Kirkpatrick (SK) model [37], which provides a conceptual framework for the understanding of the role of disorder and frustration in widely diverse systems ranging from materials to combinatorial optimization and machine learning. The SK Hamiltonian is given by

$$H_{\text{target}} = -\frac{1}{2} \sum_{i \neq j} \frac{J_{ij}}{\sqrt{N}} \sigma_i \sigma_j, \quad (9)$$

where  $\{J_{ij}\}$  is a symmetric matrix such that each matrix element  $J_{ij}$  is sampled from a gaussian distribution with mean 0 and variance 1.

Since VCA performed best in our previous examples, we use it to find ground states of the SK model for  $N = 100$  spins. Here, exact ground states energies of the SK model are calculated using the spin-glass server [30] on a total of 25 instances of disorder. To account for long-distance dependencies between spins in the SK model, we use a dilated RNN ansatz that has  $\lceil \log_2(N) \rceil = 7$  layers (see Methods Sec. VA) and set the initial temperature  $T_0 = 2$ . We compare our results with SA and SQA. For SQA, we start with an initial magnetic field  $\Gamma_0 = 2$ , while for SA we use  $T_0 = 2$ .

For an effective comparison, we first plot the residual energy per site as a function of  $N_{\text{annealing}}$  for VCA, SA and SQA (with  $P = 100$  trotter slices). Here, the SA and SQA residual energies are obtained by averaging the outcome of 50 independent annealing runs, while for VCA we average the outcome of  $10^6$  exact samples from the annealed RNN. For all methods, we take the typical average over 25 disorder instances. The results are shown in Fig. 5(a). As observed in the EA model, we note that SA and SQA produce lower residual energy solutions than VCA for small  $N_{\text{annealing}}$ , but we emphasize that VCA delivers a lower residual energy compared to SQA and SA as the total number of annealing steps increases past  $N_{\text{annealing}} \sim 10^3$ . Likewise, we observe that the rate at which the residual energy improves with increasing  $N_{\text{annealing}}$  is significantly higher for VCA in comparison to SQA and SA.

A more detailed look at the statistical behaviour of the methods at large  $N_{\text{annealing}}$  can be obtained from the residual energy histograms separately produced by each method, as shown in Fig. 5(d). The histograms contain 1000 residual energies for each of the same 25 disorder realizations. For each instance, we plot results for 1000 SA runs, 1000 samples obtained from the RNN at the end of annealing for VCA, and 10 SQA runs including contribution from each of the  $P = 100$  Trotter slices. We observe that VCA is superior to SA and SQA, as it produces a higher density of low energy configurations. This indicates that, even though VCA typically takes more annealing steps, it ultimately results in a higher chance of getting more accurate solutions to optimization problems than SA and SQA. Note that for the SK model, the SQA histogram remain quantitatively the same for 200 runs, and we report data of 10 runs only for fairness purposes compared to both SA and VCA.

We now focus on the Wishart planted ensemble (WPE), which is a class of zero-field Ising models with a first-order phase transition and tunable algorithmic hardness [38]. These problems belong to a special class of hard problem ensembles whose solutions are known a priori, which, together with the tunability of the hardness, makes the WPE model an ideal tool to benchmark heuristic algorithms for optimization problems. The Hamiltonian of the WPE model is defined as

$$H_{\text{target}} = -\frac{1}{2} \sum_{i \neq j} J_{ij}^\alpha \sigma_i \sigma_j. \quad (10)$$

Here  $J_{ij}^\alpha$  is a symmetric matrix satisfying

$$J^\alpha = \tilde{J}^\alpha - \text{diag}(\tilde{J})$$

and

$$\tilde{J}^\alpha = -\frac{1}{N} W_\alpha W_\alpha^T.$$

The term  $W_\alpha$  is an  $N \times \lfloor \alpha N \rfloor$  random matrix satisfying  $W_\alpha t_{\text{ferro}} = 0$  where  $t_{\text{ferro}} = (+1, +1, \dots, +1)$  is the ferromagnetic state (see Ref. [38] for details about the generation of  $W_\alpha$ ). The ground state of the WPE model is known (i.e., it is planted) and corresponds to the ferromagnetic states  $\pm t_{\text{ferro}}$ . Interestingly,  $\alpha$  is a tunable parameter of hardness, where for  $\alpha < 1$  this model displays a first-order transition, such that near zero temperature the paramagnetic states are meta-stable solutions [38]. This feature makes this model hard to solve with any annealing method, as the paramagnetic states are numerous compared to the two ferromagnetic states and hence act as a trap for a typical annealing method. We benchmark the three methods (SA, SQA and VCA) for  $N = 32$  and  $\alpha \in \{0.25, 0.5\}$ .

We consider 25 instances of the couplings  $\{J_{ij}^\alpha\}$  and attempt to solve the model with VCA implemented using a dilated RNN ansatz with  $\lceil \log_2(N) \rceil = 5$  layers and an initial temperature  $T_0 = 1$ . For SQA ( $P = 100$  trotter

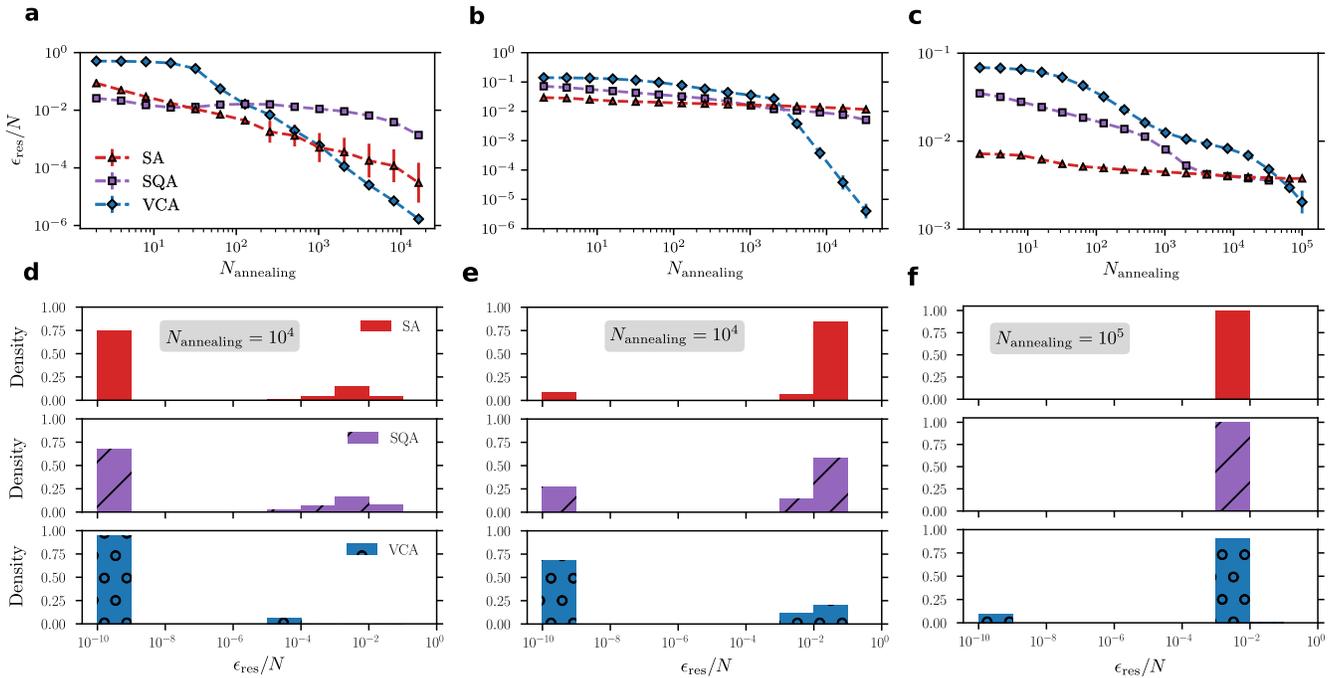


Figure 5. Benchmarking SA, SQA ( $P = 100$  trotter slices) and VCA on the Sherrington-Kirkpatrick (SK) model and the Wishart planted ensemble (WPE). Panels (a),(b), and (c) display the residual energy per site as a function of  $N_{\text{annealing}}$ . (a) The SK model with  $N = 100$  spins. (b) WPE with  $N = 32$  spins and  $\alpha = 0.5$ . (c) WPE with  $N = 32$  spins and  $\alpha = 0.25$ . Panels (d), (e) and (f) display the residual energy histogram for each of the different techniques and models in panels (a),(b), and (c), respectively. The histograms use 25000 data points for each method. Note that we choose a minimum threshold of  $10^{-10}$  for  $\epsilon_{\text{res}}/N$ , which is within our numerical accuracy.

slices), we use an initial magnetic field  $\Gamma_0 = 1$ , and for SA we start with  $T_0 = 1$ .

We first plot the scaling of residual energies per site  $\epsilon_{\text{res}}/N$  as shown in Figs. 5(b) and (c). Here we note that VCA is superior to SA and SQA for  $\alpha = 0.5$  as demonstrated in Fig. 5(b). More specifically, VCA is about three orders of magnitude more accurate than SQA and SA for a large number of annealing steps. In the case of  $\alpha = 0.25$  in Fig. 5(c), VCA is competitive where it achieves a similar performance compared to SA and SQA on average for a large number of annealing steps. We also represent the residual energies in a histogram form. We observe that for  $\alpha = 0.5$  in Fig. 5(e), VCA achieves a higher density toward low residual energies  $\epsilon_{\text{res}}/N \sim 10^{-9}$ - $10^{-10}$  compared to SA and SQA. For  $\alpha = 0.25$  in Fig. 5(f), VCA leads to a non-negligible density at very low residual energies as opposed to SA and SQA, whose solutions display residual energies orders of magnitude higher. Finally, our WPE simulations support the observation that VCA tends to improve the quality of solutions faster than SQA and SA for a large number of annealing steps.

#### IV. CONCLUSIONS AND OUTLOOK

In conclusion, we have introduced a strategy to combat the slow sampling dynamics encountered by simulated annealing when an optimization landscape is rough or glassy. Based on annealing the variational parameters of a generalized target distribution, our scheme — which we dub *variational neural annealing* — takes advantage of the power of modern autoregressive models, which can be exactly sampled without slow dynamics even when a rough landscape is encountered. We implement variational neural annealing parameterized by a recurrent neural network, and compare its performance to conventional simulated annealing on prototypical spin glass Hamiltonians known to have landscapes of varying roughness. We find that variational neural annealing produces accurate solutions to all of the optimization problems considered, including spin glass Hamiltonians where our techniques typically reach solutions orders of magnitude more accurate on average than conventional simulated annealing in the limit of a large number of annealing steps.

We emphasize that several hyperparameters, model, hardware, and variational objective function choices can be explored and may improve our methodologies. We have utilized a simple annealing schedule in our protocols and highlight that reinforcement learning can be used to improve it [39]. A critical insight gleaned from our exper-

iments is that certain neural network architectures were more efficient on specific Hamiltonians. Thus, a natural direction is to study the intimate relation between the model architecture and the problem Hamiltonian, where we envision that symmetries and domain knowledge would guide the design of models and algorithms.

As we witness the unfolding of a new age for optimization powered by deep learning [40], we anticipate a rapid adoption of machine learning techniques in the space of combinatorial optimization, as well as anticipate domain-specific applications of our ideas in diverse technological and scientific areas related to physics, biology, health care, economy, transportation, manufacturing, supply chain, hardware design, computing and information technology, among others.

## V. METHODS

### A. Recurrent Neural Network Ansätze

Recurrent neural networks model complex probability distributions  $p$  by taking advantage of the chain rule

$$p(\boldsymbol{\sigma}) = p(\sigma_1)p(\sigma_2|\sigma_1)\cdots p(\sigma_N|\sigma_{N-1},\dots,\sigma_2,\sigma_1), \quad (11)$$

where specifying every conditional probability  $p(\sigma_i|\sigma_{<i})$  provides a full characterization of the joint distribution  $p(\boldsymbol{\sigma})$ . Here,  $\{\sigma_n\}$  are  $N$  binary variables such that  $\sigma_n = 0$  corresponds to a spin down while  $\sigma_n = 1$  corresponds to a spin up. RNNs consist of elementary cells that parameterize the conditional probabilities. In their original form, “vanilla” RNN cells [41] compute a new “hidden state”  $\mathbf{h}_n$  with dimension  $d_h$ , for each site  $n$ , following the relation

$$\mathbf{h}_n = F(W[\mathbf{h}_{n-1}; \boldsymbol{\sigma}_{n-1}] + \mathbf{b}), \quad (12)$$

where  $[\mathbf{h}_{n-1}; \boldsymbol{\sigma}_{n-1}]$  is vector concatenation of  $\mathbf{h}_{n-1}$  and a one-hot encoding  $\boldsymbol{\sigma}_{n-1}$  of the binary variable  $\sigma_{n-1}$  [20]. The function  $F$  is a non-linear activation function. From this recursion relation, it is clear that the hidden state  $\mathbf{h}_n$  encodes information about the previous spins  $\sigma_{n'<n}$ . Hence, the hidden state  $\mathbf{h}_n$  provides a simple strategy to model the conditional probability  $p_\lambda(\sigma_n|\sigma_{<n})$  as

$$p_\lambda(\sigma_n|\sigma_{<n}) = \text{Softmax}(U\mathbf{h}_n + \mathbf{c}) \cdot \boldsymbol{\sigma}_n, \quad (13)$$

where  $\cdot$  denotes the dot product operation (see Fig. 6(a)). The set of all variational parameters of the model  $\boldsymbol{\lambda}$  corresponds to  $U, W, \mathbf{b}, \mathbf{c}$ , and

$$\text{Softmax}(\mathbf{v})_n = \frac{\exp(v_n)}{\sum_i \exp(v_i)}.$$

The joint probability distribution  $p_\lambda(\boldsymbol{\sigma})$  is given by

$$p_\lambda(\boldsymbol{\sigma}) = p_\lambda(\sigma_1)p_\lambda(\sigma_2|\sigma_1)\cdots p_\lambda(\sigma_N|\sigma_{<N}). \quad (14)$$

Since the outputs of the Softmax activation function sum to one, each conditional probability  $p_\lambda(\sigma_i|\sigma_{<i})$  is normalized, and hence  $p_\lambda(\boldsymbol{\sigma})$  is also normalized.

For disordered systems, it is natural to forgo the common practice of weight sharing [41] of  $W, U, \mathbf{b}$  and  $\mathbf{c}$  in Eqs. (12), (13) and use an extended set of site-dependent variational parameters  $\boldsymbol{\lambda}$  comprised of  $\{W_n\}_{n=1}^N$  and  $\{U_n\}_{n=1}^N$  and biases  $\{\mathbf{b}_n\}_{n=1}^N, \{\mathbf{c}_n\}_{n=1}^N$ . The recursion relation and the Softmax layer are modified to

$$\mathbf{h}_n = F(W_n[\mathbf{h}_{n-1}; \boldsymbol{\sigma}_{n-1}] + \mathbf{b}_n), \quad (15)$$

and

$$p_\lambda(\sigma_n|\sigma_{<n}) = \text{Softmax}(U_n\mathbf{h}_n + \mathbf{c}_n) \cdot \boldsymbol{\sigma}_n, \quad (16)$$

respectively. Note that the advantage of not using weight sharing for disordered systems is further demonstrated in Appendix. D.

We also consider a tensorized version of vanilla RNNs which replaces the concatenation operation in Eq. (15) with the operation [42]

$$\mathbf{h}_n = F(\boldsymbol{\sigma}_{n-1}^\top T_n \mathbf{h}_{n-1} + \mathbf{b}_n), \quad (17)$$

where  $\boldsymbol{\sigma}^\top$  is the transpose of  $\boldsymbol{\sigma}$ , and the variational parameters  $\boldsymbol{\lambda}$  are  $\{T_n\}_{n=1}^N, \{U_n\}_{n=1}^N, \{\mathbf{b}_n\}_{n=1}^N$  and  $\{\mathbf{c}_n\}_{n=1}^N$ . This form of tensorized RNN increases the expressiveness of our ansatz as illustrated in Appendix. D.

For two-dimensional systems, we make use of a 2D-dimensional extension of the recursion relation in vanilla RNNs [20]

$$\mathbf{h}_{i,j} = F\left(W_{i,j}^{(h)}[\mathbf{h}_{i-1,j}; \boldsymbol{\sigma}_{i-1,j}] + W_{i,j}^{(v)}[\mathbf{h}_{i,j-1}; \boldsymbol{\sigma}_{i,j-1}] + \mathbf{b}_{i,j}\right). \quad (18)$$

To enhance the expressive power of the model, we promote the recursion relation to a tensorized form

$$\mathbf{h}_{i,j} = F([\boldsymbol{\sigma}_{i-1,j}; \boldsymbol{\sigma}_{i,j-1}]T_{i,j}[\mathbf{h}_{i-1,j}; \mathbf{h}_{i,j-1}] + \mathbf{b}_{i,j}). \quad (19)$$

Here,  $T_{i,j}$  are site-dependent weight tensors that have dimension  $4 \times 2d_h \times d_h$ . We also note that the coordinates  $(i-1, j)$  and  $(i, j-1)$  are path-dependent, and are given by the zigzag path, illustrated by the black arrows in Fig. 6(b). Moreover, to sample configurations from the 2D tensorized RNNs, we use the same zigzag path as illustrated by the red dashed arrows in Fig. 6(b).

For models such as the Sherrington-Kirkpatrick model and the Wishart planted ensemble, every spin interacts with each other. To account for the long-distance nature of the correlations induced by these interactions, we use dilated RNNs [43], which are known to alleviate the vanishing gradient problem [44]. Dilated RNNs are multi-layered RNNs that use dilated connections between spins to model long-term dependencies [45], as illustrated in Fig. 6(c). At each layer  $1 \leq l \leq L$ , the hidden state is computed as

$$\mathbf{h}_n^{(l)} = F(W_n^{(l)}[\mathbf{h}_{\max(0, n-2^{l-1})}^{(l-1)}; \mathbf{h}_n^{(l-1)}] + \mathbf{b}_n^{(l)}).$$

Here  $\mathbf{h}_n^{(0)} = \boldsymbol{\sigma}_{n-1}$  and the conditional probability is given by

$$p_\lambda(\sigma_n|\sigma_{<n}) = \text{Softmax}(U_n\mathbf{h}_n^{(L)} + \mathbf{c}_n) \cdot \boldsymbol{\sigma}_n.$$

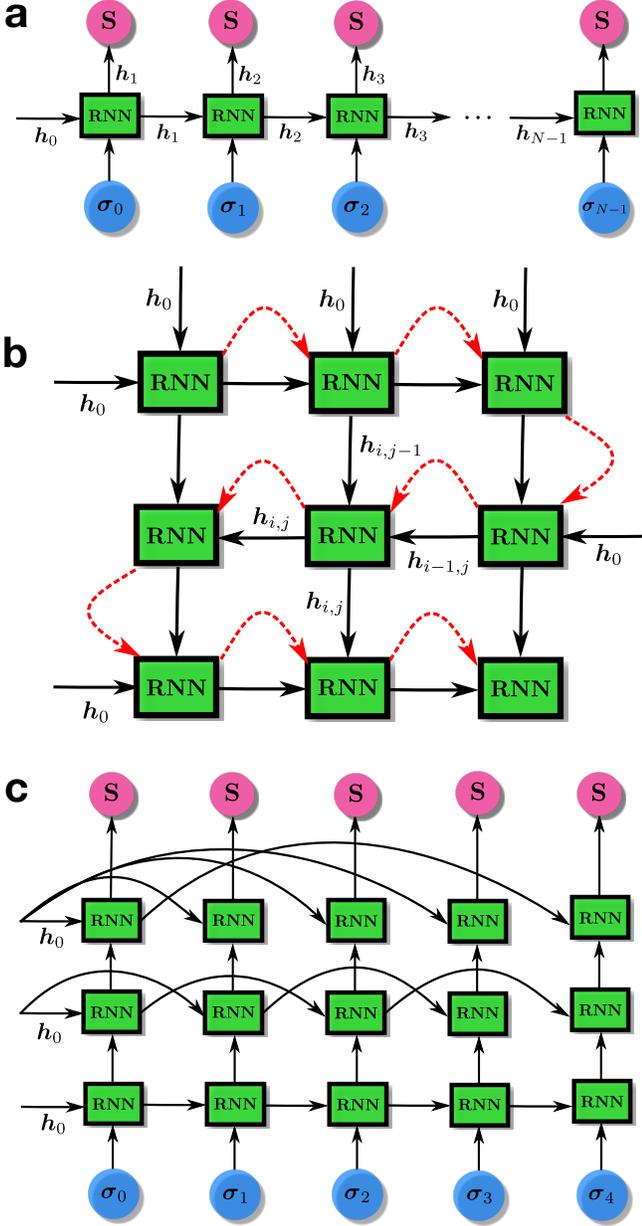


Figure 6. (a) An illustration of a 1D RNN: at each site  $n$ , the RNN cell denoted by the green box, receives a hidden state  $\mathbf{h}_{n-1}$  and the one-hot spin vector  $\boldsymbol{\sigma}_{n-1}$ , to generate a new hidden state  $\mathbf{h}_n$  that is fed into a Softmax layer (denoted by a magenta circle). (b) A graphical illustration of a 2D RNN. Each RNN cell receives two hidden states  $\mathbf{h}_{i,j-1}$  and  $\mathbf{h}_{i-1,j}$ , as well as two input vectors  $\boldsymbol{\sigma}_{i,j-1}$  and  $\boldsymbol{\sigma}_{i-1,j}$  (not shown) as illustrated by the black arrows. The red arrows correspond to the zigzag path we use for 2D autoregressive sampling. The initial memory state  $\mathbf{h}_0$  of the RNN and the initial inputs  $\boldsymbol{\sigma}_0$  (not shown) are null vectors. (c) An illustration of a dilated RNN, where the distance between each two RNN cells grows exponentially with depth to account for long-term dependencies. We choose depth  $L = \lceil \log_2(N) \rceil$  where  $N$  is the number of spins.

In our work, we choose the size of the hidden states  $\mathbf{h}_n^{(l)}$ , where  $l > 0$ , as constant and equal to  $d_h$ . We also use a number of layers  $L = \lceil \log_2(N) \rceil$ , where  $N$  is the number of spins and  $\lceil \dots \rceil$  is the ceiling function. This means that two spins are connected with a path whose length is bounded by  $\mathcal{O}(\log_2(N))$ , which follows the spirit of the multi-scale renormalization ansatz [46]. For more details on the advantage of dilated RNNs over tensorized RNNs see Appendix. D.

We finally note that for all the RNN architectures in our work, we found accurate results using the exponential linear unit (ELU) activation function, defined as:

$$\text{ELU}(x) = \begin{cases} x, & \text{if } x \geq 0, \\ \exp(x) - 1, & \text{if } x < 0. \end{cases}$$

## B. Minimizing the variational free energy

To implement the variational classical annealing algorithm, we use the variational free energy

$$F_\lambda(T) = \langle H_{\text{target}} \rangle_\lambda - T S_{\text{classical}}(p_\lambda), \quad (20)$$

where the target Hamiltonian  $H_{\text{target}}$  encodes the optimization problem and  $T$  is the temperature. Moreover,  $S_{\text{classical}}$  is the entropy of the distribution  $p_\lambda$ . To estimate  $F_\lambda(T)$  we take  $N_s$  exact samples  $\boldsymbol{\sigma}^{(i)} \sim p_\lambda$  ( $i = 1, \dots, N_s$ ) drawn from the RNN and evaluate

$$F_\lambda(T) \approx \frac{1}{N_s} \sum_{i=1}^{N_s} F_{\text{loc}}(\boldsymbol{\sigma}^{(i)}),$$

where the local free energy is  $F_{\text{loc}}(\boldsymbol{\sigma}) = H_{\text{target}}(\boldsymbol{\sigma}) + T \log(p_\lambda(\boldsymbol{\sigma}))$  [18]. Similarly, the gradients are given by

$$\begin{aligned} \partial_\lambda F_\lambda(T) \approx & \frac{1}{N_s} \sum_{i=1}^{N_s} \partial_\lambda \log(p_\lambda(\boldsymbol{\sigma}^{(i)})) \\ & \times (F_{\text{loc}}(\boldsymbol{\sigma}^{(i)}) - F_\lambda(T)), \end{aligned}$$

where we subtract  $F_\lambda(T)$  in order to reduce noise in the gradients [18, 20]. We note that this variational scheme exhibits a zero-variance principle, namely that the local free energy variance per spin

$$\sigma_F^2 \equiv \frac{\text{var}(\{F_{\text{loc}}(\boldsymbol{\sigma})\})}{N}, \quad (21)$$

becomes zero when  $p_\lambda$  matches the Boltzmann distribution, provided that mode collapse is avoided [18].

The gradient updates are implemented using the Adam optimizer [47]. Furthermore, the computational complexity of VCA for one gradient descent step is  $\mathcal{O}(N_s \times N \times d_h^2)$  for 1D RNNs and 2D RNNs (both vanilla and tensorized versions) and  $\mathcal{O}(N_s \times N \log(N) \times d_h^2)$  for dilated RNNs. Consequently, VCA has lower computational cost than VQA, which is implemented using VMC (see Methods Sec. VC).

Finally, we note that in our implementations no training steps are performed at the end of annealing for both VCA and VQA.

### C. Variational Monte Carlo

The main goal of Variational Monte Carlo is to approximate the ground state of a Hamiltonian  $\hat{H}$  through the iterative optimization of an ansatz wave function  $|\Psi_\lambda\rangle$ . The VMC objective function is given by

$$E \equiv \frac{\langle \Psi_\lambda | \hat{H} | \Psi_\lambda \rangle}{\langle \Psi_\lambda | \Psi_\lambda \rangle}.$$

We note that an important class of stoquastic many-body Hamiltonians has ground states  $|\Psi\rangle$  with strictly real and positive amplitudes in the standard product spin basis [48]. These ground states can be written down in terms of probability distributions,

$$|\Psi\rangle = \sum_{\boldsymbol{\sigma}} \Psi(\boldsymbol{\sigma}) |\boldsymbol{\sigma}\rangle = \sum_{\boldsymbol{\sigma}} \sqrt{P(\boldsymbol{\sigma})} |\boldsymbol{\sigma}\rangle. \quad (22)$$

To approximate this family of states, we use an RNN wave function, namely  $\Psi_\lambda(\boldsymbol{\sigma}) = \sqrt{p_\lambda(\boldsymbol{\sigma})}$ . Extensions to complex-valued RNN wave functions are defined in Ref. [20], and results on their ability to simulate variational quantum annealing of non-stoquastic Hamiltonians [49] will be reported elsewhere [50]. These families of RNN states are normalized by construction (i.e.,  $\langle \Psi_\lambda | \Psi_\lambda \rangle = 1$ ) and allow for accurate estimates of the energy expectation value. By taking  $N_s$  exact samples  $\boldsymbol{\sigma}^{(i)} \sim p_\lambda$  ( $i = 1, \dots, N_s$ ), it follows that

$$E \approx \frac{1}{N_s} \sum_{i=1}^{N_s} E_{\text{loc}}(\boldsymbol{\sigma}^{(i)}).$$

The local energy is given by

$$E_{\text{loc}}(\boldsymbol{\sigma}) = \sum_{\boldsymbol{\sigma}'} H_{\boldsymbol{\sigma}\boldsymbol{\sigma}'} \frac{\Psi_\lambda(\boldsymbol{\sigma}')}{\Psi_\lambda(\boldsymbol{\sigma})}, \quad (23)$$

where the sum over  $\boldsymbol{\sigma}'$  is tractable when the Hamiltonian  $\hat{H}$  is local. Similarly, we can also estimate the energy gradients as

$$\partial_\lambda E = \frac{2}{N_s} \sum_{i=1}^{N_s} \partial_\lambda \log(\Psi_\lambda(\boldsymbol{\sigma}^{(i)})) (E_{\text{loc}}(\boldsymbol{\sigma}^{(i)}) - E).$$

Here, we can subtract the term  $E$  in order to reduce noise in the stochastic estimation of our gradients without introducing a bias [20, 51]. In fact, when the ansatz is close to an eigenstate of  $\hat{H}$ , then  $E_{\text{loc}}(\boldsymbol{\sigma}) \approx E$ , which means that the variance of gradients  $\text{Var}(\partial_{\lambda_j} E) \approx 0$  for each variational parameter  $\lambda_j$ . We note that this is similar in spirit to the control variate methods in Monte Carlo and to the baseline methods in reinforcement learning [51].

Similarly to the minimization scheme of the variational free energy in Methods Sec. VB, VMC also exhibits a zero-variance principle, where the energy variance per spin

$$\sigma^2 \equiv \frac{\text{var}(\{E_{\text{loc}}(\boldsymbol{\sigma})\})}{N}, \quad (24)$$

becomes zero when  $|\Psi_\lambda\rangle$  matches an excited state of  $\hat{H}$ , which thanks to the minimization of the variational energy  $E$  is likely to be the ground state  $|\Psi_G\rangle$ .

The gradients  $\partial_\lambda \log(\Psi_\lambda(\boldsymbol{\sigma}))$  are numerically computed using automatic differentiation [52]. We use the Adam optimizer to perform gradient descent updates, with a learning rate  $\eta$ , to optimize the variational parameters  $\boldsymbol{\lambda}$  of the RNN wave function. We note that in the presence of  $\mathcal{O}(N)$  non-diagonal elements in a Hamiltonian  $\hat{H}$ , the local energies  $E_{\text{loc}}(\boldsymbol{\sigma})$  have  $\mathcal{O}(N)$  terms (see Eq. (23)). Thus, the computational complexity of one gradient descent step is  $\mathcal{O}(N_s \times N^2 \times d_h^2)$  for 1D RNNs and 2D RNNs (both vanilla and tensorized versions).

### D. Simulated Quantum Annealing and Simulated Annealing

Simulated Quantum Annealing is a standard quantum-inspired classical technique that has traditionally been used to benchmark the behavior of quantum annealers [24]. It is usually implemented via the path-integral Monte Carlo method [11], a QMC method that simulates equilibrium properties of quantum systems at finite temperature. To illustrate this method, consider a  $D$ -dimensional time-dependent quantum Hamiltonian

$$\hat{H}(t) = - \sum_{i,j} J_{ij} \hat{\sigma}_i^z \hat{\sigma}_j^z - \Gamma(t) \sum_{i=1}^N \hat{\sigma}_i^x,$$

where  $\Gamma(t) = \Gamma_0(1-t)$  controls the strength of the quantum annealing dynamics at a time  $t \in [0, 1]$ . By applying the Suzuki-Trotter formula to the partition function of the quantum system,

$$Z = \text{Tr} \exp\{-\beta \hat{H}(t)\}, \quad (25)$$

with the inverse temperature  $\beta = \frac{1}{T}$ , we can map the  $D$ -dimensional quantum Hamiltonian onto a  $(D+1)$  classical system consisting of  $P$  coupled replicas (Trotter slices) of the original system

$$H_{D+1}(t) = - \sum_{k=1}^P \left( \sum_{i,j} J_{ij} \sigma_i^k \sigma_j^k + J_\perp(t) \sum_{i=1}^N \sigma_i^k \sigma_i^{k+1} \right), \quad (26)$$

where  $\sigma_i^k$  is the classical spin at site  $i$  and replica  $k$ . The term  $J_\perp(t)$  corresponds to uniform coupling between  $\sigma_i^k$  and  $\sigma_i^{k+1}$  for each site  $i$ , such that

$$J_\perp(t) = -\frac{PT}{2} \ln \left( \tanh \left( \frac{\Gamma(t)}{PT} \right) \right).$$

We note that periodic boundary conditions  $\sigma^{P+1} \equiv \sigma^1$  arise because of the trace in Eq. (25).

Interestingly, we can approximate  $Z$  with an effective partition function  $Z_p$  at temperature  $PT$  given by [35]:

$$Z_p \propto \text{Tr} \exp \left\{ -\frac{H_{D+1}(t)}{PT} \right\},$$

which can now be simulated with a standard Metropolis-Hastings Monte Carlo algorithm. A key element to this algorithm is the energy difference induced by a single spin flip at site  $\sigma_i^k$ , which is equal to

$$\Delta_i E_{\text{local}} = 2 \sum_j J_{ij} \sigma_i^k \sigma_j^k + 2J_{\perp}(t) (\sigma_i^{k-1} \sigma_i^k + \sigma_i^k \sigma_i^{k+1}).$$

Here, the second term encodes the quantum dynamics. In our simulations we consider single spin flip (local) moves applied to all sites in all slices. We can also perform a global move [35], which means flipping a spin at location  $i$  in every slice  $k$ . Clearly this has no impact on the term dependent on  $J_{\perp}$ , because it contains only terms quadratic in the flipped spin, so that

$$\Delta_i E_{\text{global}} = 2 \sum_{k=1}^P \sum_j J_{ij} \sigma_i^k \sigma_j^k.$$

In summary, a single Monte Carlo step (MCS) consists of first performing a single local move on all sites in each  $k$ -th slice and on all slices, followed by a global move for all sites. For the SK model and the WPE model studied in this paper, we use  $P = 100$ , whereas for the EA model we use  $P = 20$  similarly to Ref. [11]. Before starting the quantum annealing schedule, we first thermalize the system by performing SA [35] from a temperature  $T_0 = 3$  to a final temperature  $1/P$  (so that  $PT = 1$ ). This is done in 60 steps, where at each temperature we perform 100 Metropolis moves on each site. We then perform SQA using a linear schedule that decreases the field from  $\Gamma_0$  to a final value close to zero  $\Gamma(t=1) = 10^{-8}$ , where five local and global moves are performed for each value of the magnetic field  $\Gamma(t)$ , so that it is consistent with the choice of  $N_{\text{train}} = 5$  for VCA (see Sec. II and III A). Thus, the number of MCS is equal to five times the number of annealing steps.

For the standalone SA, we decrease the temperature from  $T_0$  to  $T(t=1) = 10^{-8}$ . Here, a single MCS consists of a Monte Carlo sweep, i.e., attempting a spin-flip for all sites. For each thermal annealing step, we perform five MCS, and hence similar to SQA, the number of MCS is equal to five times the number of annealing steps. Furthermore, we do a warm-up step for SA, by performing  $N_{\text{warmup}}$  MCS to equilibrate the Markov Chain at the initial temperature  $T_0$  and to provide a consistent choice with VCA (see Sec. II).

## ACKNOWLEDGMENTS

We acknowledge Jack Raymond for suggesting to use the Wishart Planted Ensemble as a benchmark for our variational annealing setup. We also thank Christopher Roth, Cunlu Zhou, Martin Ganahl and Giuseppe Santoro for fruitful discussions. We are also grateful to Lauren Hayward for providing her plotting code to produce our figures using Matplotlib library. Our RNN implementation is based on Tensorflow and NumPy. We acknowledge support from the Natural Sciences and Engineering Research Council (NSERC), a Canada Research Chair, the Shared Hierarchical Academic Research Computing Network (SHARCNET), Compute Canada, Google Quantum Research Award, and the Canadian Institute for Advanced Research (CIFAR) AI chair program. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute [www.vectorinstitute.ai/#partners](http://www.vectorinstitute.ai/#partners). Research at Perimeter Institute is supported in part by the Government of Canada through the Department of Innovation, Science and Economic Development Canada and by the Province of Ontario through the Ministry of Economic Development, Job Creation and Trade.

## Appendix A: Numerical proof of principle of adiabaticity

As demonstrated in Sec. III, we have shown that both VQA and VCA are effective at finding the classical ground state of disordered spin chains. Here, we further illustrate the adiabaticity of both VQA and VCA. First, we perform VQA on the uniform ferromagnetic Ising chain (i.e.,  $J_{i,i+1} = 1$ ) with  $N = 20$  spins and open boundary conditions with an initial transverse field  $\Gamma_0 = 2$ . Here, we use a tensorized RNN wave function with weight sharing across sites of the chain. We also choose  $N_{\text{annealing}} = 1024$ . In Fig. 7(a), we show that the variational energy tracks the exact ground energy throughout the annealing process with high accuracy. We also observe that optimizing an RNN wave function from scratch, i.e., randomly reinitializing the parameters of the model at each new value of the transverse magnetic field is not optimal. This observation underlines the importance of transferring the parameters of our wave function ansatz after each annealing step. Furthermore, in Fig. 7(b) we illustrate that the RNN wave function's residual energy is much lower compared to the gap throughout the annealing process, which shows that VQA remains adiabatic for a large number of annealing steps.

Similarly, in Fig. 7(c) we perform VCA with an initial temperature  $T_0 = 2$  on the same model, the same system size, the same ansatz, and the same number of annealing steps. We see an excellent agreement between the RNN wave function free energy and the exact free energy, high-

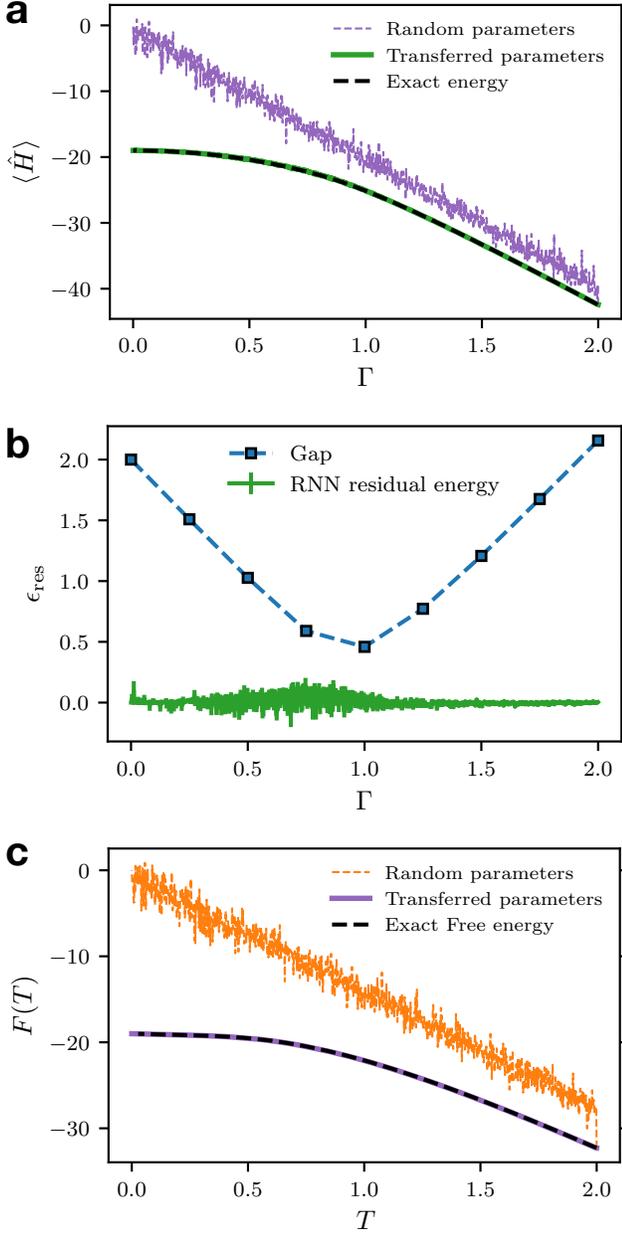


Figure 7. Numerical evidence of adiabaticity on the uniform Ising chain with  $N = 20$  spins for VQA in panels (a) and (b) and VCA in panel (c). (a) Variational energy of RNN wave function against the transverse magnetic field  $\Gamma$ , with  $\lambda$  initialized using the parameters optimized in the previous annealing step (transferred parameters, green curve) and with random parameter reinitialization (random parameters, purple curve). These strategies are compared with the exact energy obtained from exact diagonalization (dashed black line). (b) Residual energy of the RNN wave function vs the transverse field  $\Gamma$ . Throughout annealing with VQA, the residual energy is always much smaller than the gap within error bars. (c) Variational free energy vs temperature  $T$  for a VCA run with  $\lambda$  initialized using the parameters optimized in the previous annealing step (transferred parameters, purple line) and with random reinitialization (random parameters, orange curve).

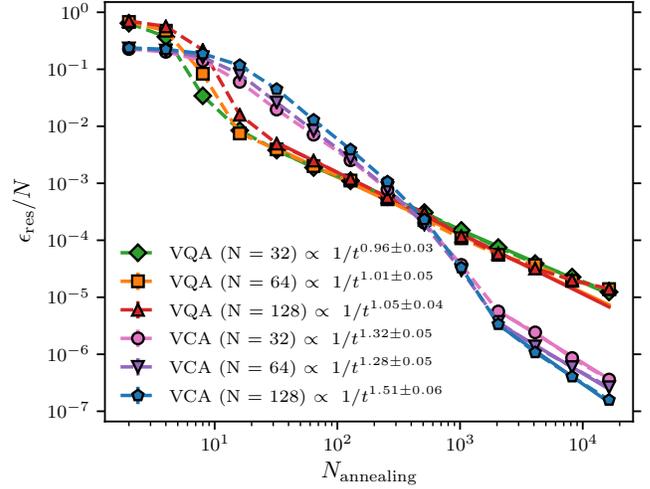


Figure 8. Variational annealing on random Ising chains, where we represent the residual energy per site  $\epsilon_{\text{res}}/N$  vs  $N_{\text{annealing}}$  for both VQA and VCA. The system sizes are  $N = 32, 64, 128$  and we use random discrete couplings  $J_{i,i+1} \in \{-1, 1\}$ .

lighting once again the adiabaticity of our emulation of classical annealing, as well as the importance of transferring the parameters of our ansatz after each annealing step. Taken all together, the results in Fig. 7 support the notion that VQA and VCA evolutions can be adiabatic.

In Fig. 8 we report the residual energies per site against the number of annealing steps  $N_{\text{annealing}}$ . Here, we consider  $J_{i,i+1}$  uniformly sampled from the discrete set  $\{-1, +1\}$ , where the ground state configuration is disordered and the ground state energy is given by  $E_G = -\sum_{i=1}^{N-1} |J_{i,i+1}| = -(N-1)$ . The decay exponents for VCA are in the interval 1.3 – 1.6 and the VQA exponent are approximately 1. These exponents also suggest an asymptotic speed-up compared to SA and coherent quantum annealing, where the residual energies follow a logarithmic law [29, 53–55]. The latter confirms the robustness of the observations in Fig. 3.

## Appendix B: The variational adiabatic theorem

In this section, we derive a sufficient condition for the number of gradient descent steps needed to maintain the variational ansatz close to the instantaneous ground state throughout the VQA simulation. First, consider a variational wave function  $|\Psi_\lambda\rangle$  and the following the time-dependent Hamiltonian:

$$\hat{H}(t) = \hat{H}_{\text{target}} + f(t)\hat{H}_D,$$

The goal is to find the ground state of the target Hamiltonian  $\hat{H}_{\text{target}}$  by introducing quantum fluctuations through a driving Hamiltonian  $\hat{H}_D$ , where  $\hat{H}_D \gg \hat{H}_{\text{target}}$ . Here  $f(t)$  is a decreasing schedule function such that  $f(0) = 1$ ,  $f(1) = 0$  and  $t \in [0, 1]$ .

Let  $E(\boldsymbol{\lambda}, t) = \langle \Psi_{\boldsymbol{\lambda}} | \hat{H}(t) | \Psi_{\boldsymbol{\lambda}} \rangle$ , and  $E_G(t), E_E(t)$  the instantaneous ground/excited state energy of the Hamiltonian  $\hat{H}(t)$ , respectively. The instantaneous energy gap is defined as  $g(t) \equiv E_E(t) - E_G(t)$ .

To simplify our discussion, we consider the case of a target Hamiltonian that has a non-degenerate ground state. Here, we decompose the variational wave function as:

$$|\Psi_{\boldsymbol{\lambda}}\rangle = (1 - a(t))^{\frac{1}{2}} |\Psi_G(t)\rangle + a(t)^{\frac{1}{2}} |\Psi_{\perp}(t)\rangle, \quad (\text{B1})$$

where  $|\Psi_G(t)\rangle$  is the instantaneous ground state and  $|\Psi_{\perp}(t)\rangle$  is a superposition of all the instantaneous excited states. From this decomposition, one can show that [56]:

$$a(t) \leq \frac{E(\boldsymbol{\lambda}, t) - E_G(t)}{g(t)}. \quad (\text{B2})$$

As a consequence, in order to satisfy adiabaticity, i.e.,  $|\langle \Psi_{\perp}(t) | \Psi_{\boldsymbol{\lambda}} \rangle|^2 \ll 1$  for all times  $t$ , then one should have  $a(t) < \epsilon \ll 1$  where  $\epsilon$  is a small upper bound on the overlap between the variational wave function and the excited states. This means that the success probability  $P_{\text{success}}$  of obtaining the ground state at  $t = 1$  is bounded from below by  $1 - \epsilon$ . From Eq. (B2), to satisfy  $a(t) < \epsilon$ , it is sufficient to have:

$$\epsilon_{\text{res}}(\boldsymbol{\lambda}, t) \equiv E(\boldsymbol{\lambda}, t) - E_G(t) < \epsilon g(t). \quad (\text{B3})$$

To satisfy the latter condition, we require a slightly stronger condition as follows:

$$\epsilon_{\text{res}}(\boldsymbol{\lambda}, t) < \frac{\epsilon g(t)}{2}. \quad (\text{B4})$$

In our derivation of a sufficient condition on the number of gradient descent steps to satisfy the previous requirement, we use the following set of assumptions:

- **(A1)**  $|\partial_t^k E_G(t)|, |\partial_t^k g(t)|, |\partial_t^k f(t)| \leq \mathcal{O}(\text{poly}(N))$ , for all  $0 \leq t \leq 1$  and for  $k \in \{1, 2\}$ .
- **(A2)**  $|\langle \Psi_{\boldsymbol{\lambda}} | \hat{H}_D | \Psi_{\boldsymbol{\lambda}} \rangle| \leq \mathcal{O}(\text{poly}(N))$  for all possible parameters  $\boldsymbol{\lambda}$  of the variational wave function.
- **(A3)** No anti-crossing during annealing, i.e.,  $g(t) \neq 0$ , for all  $0 \leq t \leq 1$ .
- **(A4)** The gradients  $\partial_{\boldsymbol{\lambda}} E(\boldsymbol{\lambda}, t)$  can be calculated exactly, are  $L(t)$ -Lipschitz with respect to  $\boldsymbol{\lambda}$  and  $L(t) \leq \mathcal{O}(\text{poly}(N))$  for all  $0 \leq t \leq 1$ .
- **(A5)** Local convexity, i.e., close to convergence when  $\epsilon_{\text{res}}(\boldsymbol{\lambda}, t) < \epsilon g(t)$ , the energy landscape of  $E(\boldsymbol{\lambda}, t)$  is convex with respect to  $\boldsymbol{\lambda}$ , for all  $0 < t \leq 1$ .

Note that this assumption is  $\epsilon$ -dependent.

- **(A6)** The parameters vector  $\boldsymbol{\lambda}$  is bounded by a polynomial in  $N$ . i.e.,  $\|\boldsymbol{\lambda}\| \leq \mathcal{O}(\text{poly}(N))$ , where we define “ $\|\cdot\|$ ” as the euclidean  $L_2$  norm.

- **(A7)** The variational wave function  $|\Psi_{\boldsymbol{\lambda}}\rangle$  is expressive enough, i.e.,

$$\min_{\boldsymbol{\lambda}} \epsilon_{\text{res}}(\boldsymbol{\lambda}, t) < \frac{\epsilon g(t)}{4}, \quad \forall t \in [0, 1].$$

Note that this assumption is also  $\epsilon$ -dependent.

- **(A8)** At  $t = 0$ , the energy landscape of  $E(\boldsymbol{\lambda}, t = 0)$  is globally convex with respect to  $\boldsymbol{\lambda}$ .

**Theorem** Given the assumptions **(A1)** to **(A8)**, a sufficient (but not necessary) number of gradient descent steps  $N_{\text{steps}}$  to satisfy the condition (B4) during the VQA protocol, is bounded as:

$$\mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon \min_{\{t_n\}}(g(t_n))}\right) \leq N_{\text{steps}} \leq \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon^2 \min_{\{t_n\}}(g(t_n))^2}\right),$$

where  $(t_1, t_2, t_3, \dots)$  is an increasing finite sequence of time steps, satisfying  $t_1 = 0$  and  $t_{n+1} = t_n + \delta t_n$ , where

$$\delta t_n = \mathcal{O}\left(\frac{\epsilon g(t_n)}{\text{poly}(N)}\right).$$

**Proof:** In order to satisfy the condition Eq. (B4) during the VQA protocol, we follow these steps:

- Step 1 (warm-up step): we prepare our variational wave function at the ground state at  $t = 0$  such that Eq. (B4) is verified at time  $t = 0$ .
- Step 2 (annealing step): we change time  $t$  by an infinitesimal amount  $\delta t$ , so that the condition (B3) is verified at time  $t + \delta t$ .
- Step 3 (training step): we tune the parameters of the variational wave function, using gradient descent, so that the condition (B4) is satisfied at time  $t + \delta t$ .
- Step 4: we loop over steps 2 and 3 until we arrive at  $t = 1$ , where we expect to obtain the ground state energy of the target Hamiltonian.

Let us first start with step 2 assuming that step 1 is verified. In order to satisfy the requirement of this step at time  $t$ , then  $\delta t$  has to be chosen small enough so that

$$\epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t + \delta t) < \epsilon g(t + \delta t) \quad (\text{B5})$$

is verified given that the condition (B4) is satisfied at time  $t$ . Here,  $\boldsymbol{\lambda}_t$  are the parameters of the variational wave function that satisfies the condition (B4) at time  $t$ . To get a sense of how small  $\delta t$  should be, we do a Taylor expansion, while fixing the parameters  $\boldsymbol{\lambda}_t$ , to get:

$$\begin{aligned} & \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t + \delta t) \\ &= \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) + \partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) \delta t + \mathcal{O}((\delta t)^2), \\ &< \frac{\epsilon g(t)}{2} + \partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) \delta t + \mathcal{O}((\delta t)^2), \end{aligned}$$

where we used the condition (B4) to go from the second line to the third line. Here,  $\partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) = \partial_t f(t) \langle \hat{H}_D \rangle - \partial_t E_G(t)$ . To satisfy the condition (B3) at time  $t + \delta t$ , it is enough to have the right hand side of the previous inequality to be much smaller than the gap at  $t + \delta t$ , i.e.,

$$\frac{\epsilon g(t)}{2} + \partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) \delta t + \mathcal{O}((\delta t)^2) < \epsilon g(t + \delta t).$$

By Taylor expanding the gap, we get:

$$\partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) \delta t + \mathcal{O}((\delta t)^2) < \frac{\epsilon g(t)}{2} + \epsilon \partial_t g(t) \delta t + \mathcal{O}((\delta t)^2),$$

hence, it is enough to satisfy the following condition:

$$(\partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) - \epsilon \partial_t g(t)) \delta t + \mathcal{O}((\delta t)^2) < \frac{\epsilon g(t)}{2}. \quad (\text{B6})$$

Using the Taylor-Laplace formula, one can express the Taylor remainder term  $\mathcal{O}((\delta t)^2)$  as follows:

$$\mathcal{O}((\delta t)^2) = \int_t^{t+\delta t} (\tau - t) A(\tau) d\tau,$$

where  $A(\tau) = \partial_\tau^2 \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, \tau) - \epsilon \partial_\tau^2 g(\tau) = \partial_\tau^2 f(\tau) \langle \hat{H}_D \rangle - \partial_\tau^2 E_G(\tau) - \epsilon \partial_\tau^2 g(\tau)$  and  $\tau$  is between  $t$  and  $t + \delta t$ . The last expression can be bounded as follows:

$$\mathcal{O}((\delta t)^2) \leq \int_t^{t+\delta t} (\tau - t) |A(\tau)| d\tau \leq \frac{(\delta t)^2}{2} \sup(|A|).$$

where “ $\sup(|A|)$ ” is the supremum of  $|A|$  over the interval  $[0, 1]$ . Given assumptions (A1) and (A2), then  $\sup(|A|)$  is bounded from above by a polynomial in  $N$ , hence:

$$\mathcal{O}((\delta t)^2) \leq \mathcal{O}(\text{poly}(N)) (\delta t)^2 \leq \mathcal{O}(\text{poly}(N)) \delta t,$$

where the last inequality holds since  $\delta t \leq 1$  as  $t \in [0, 1]$ , while we note that it is not necessarily tight. Furthermore, since  $(\partial_t \epsilon_{\text{res}}(\boldsymbol{\lambda}_t, t) - \epsilon \partial_t g(t))$  is also bounded from above by a polynomial in  $N$  (according to assumptions (A1) and (A2)), then in order to satisfy Eq. (B6), it is sufficient to require the following condition:

$$\mathcal{O}(\text{poly}(N)) \delta t < \frac{\epsilon g(t)}{2}.$$

Thus, it is sufficient to take:

$$\delta t = \mathcal{O}\left(\frac{\epsilon g(t)}{\text{poly}(N)}\right). \quad (\text{B7})$$

By taking account of assumption (A3),  $\delta t$  can be taken non-zero for all time steps  $t$ . As a consequence, assuming the condition (B7) is verified for a non-zero  $\delta t$  and a suitable  $\mathcal{O}(1)$  prefactor, then the condition (B5) is also verified.

We can now move to step 3. Here, we apply a number of gradient descent steps  $N_{\text{train}}(t)$  to find a new set of parameters  $\boldsymbol{\lambda}_{t+\delta t}$  such that:

$$\epsilon_{\text{res}}(\boldsymbol{\lambda}_{t+\delta t}, t+\delta t) = E(\boldsymbol{\lambda}_{t+\delta t}, t+\delta t) - E_G(t+\delta t) < \frac{\epsilon g(t + \delta t)}{2}, \quad (\text{B8})$$

To estimate the scaling of the number of gradient descent steps  $N_{\text{train}}(t)$  needed to satisfy (B8), we make use of assumptions (A4) and (A5). The assumption (A5) is reasonable providing that the variational energy  $E(\boldsymbol{\lambda}_t, t + \delta t)$  is very close to the ground state energy  $E_G(t + \delta t)$ , as given by Eq. (B5). Using the above assumptions and assuming that the learning rate  $\eta(t) = 1/L(t)$ , we can use a well-known result in convex optimization [57] (see Sec. 2.1.5), which states the following inequality:

$$E(\tilde{\boldsymbol{\lambda}}_t, t + \delta t) - \min_{\boldsymbol{\lambda}} E(\boldsymbol{\lambda}, t + \delta t) \leq \frac{2L(t) \|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t+\delta t}^*\|^2}{N_{\text{train}}(t) + 4}.$$

Here,  $\tilde{\boldsymbol{\lambda}}_t$  are the new variational parameters obtained after applying  $N_{\text{train}}(t + \delta t)$  gradient descent steps starting from  $\boldsymbol{\lambda}_t$ . Furthermore,  $\boldsymbol{\lambda}_{t+\delta t}^*$  are the optimal parameters such that:

$$E(\boldsymbol{\lambda}_{t+\delta t}^*, t + \delta t) = \min_{\boldsymbol{\lambda}} E(\boldsymbol{\lambda}, t + \delta t).$$

Since the Lipschitz constant  $L(t) \leq \mathcal{O}(\text{poly}(N))$  (assumption (A4)) and  $\|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t+\delta t}^*\|^2 \leq \mathcal{O}(\text{poly}(N))$  (assumption (A6)), one can take

$$N_{\text{train}}(t + \delta t) = \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon g(t + \delta t)}\right), \quad (\text{B9})$$

with a suitable  $\mathcal{O}(1)$  prefactor, so that:

$$E(\tilde{\boldsymbol{\lambda}}_t, t + \delta t) - \min_{\boldsymbol{\lambda}} E(\boldsymbol{\lambda}, t + \delta t) < \frac{\epsilon g(t + \delta t)}{4}.$$

Moreover, by assuming that the variational wave function is expressive enough (assumption (A7)), i.e.,

$$\min_{\boldsymbol{\lambda}} E(\boldsymbol{\lambda}, t + \delta t) - E_G(t + \delta t) < \frac{\epsilon g(t + \delta t)}{4},$$

we can then deduce, by taking  $\boldsymbol{\lambda}_{t+\delta t} \equiv \tilde{\boldsymbol{\lambda}}_t$  and summing the two previous inequalities, that:

$$E(\boldsymbol{\lambda}_{t+\delta t}, t + \delta t) - E_G(t + \delta t) < \frac{\epsilon g(t + \delta t)}{2}.$$

Let us recall that in step 1, we have to initially prepare the variational ansatz to satisfy condition (B4) at  $t = 0$ . In fact, we can take advantage of the assumption (A4), where the gradients are  $L(0)$ -Lipschitz with  $L(0) \leq \mathcal{O}(\text{poly}(N))$ . We can also use the convexity assumption (A8), and we can show that a sufficient number of gradient descent steps to satisfy condition (B4) at  $t = 0$  is estimated as:

$$N_{\text{warmup}} \equiv N_{\text{train}}(0) = \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon g(0)}\right).$$

The latter can be obtained in a similar way as in Eq. (B9).

In conclusion, the total number of gradient steps  $N_{\text{steps}}$  to evolve the Hamiltonian  $\hat{H}(0)$  to the target Hamiltonian  $\hat{H}(1)$ , while verifying the condition (B4) is given by:

$$N_{\text{steps}} = \sum_{n=1}^{N_{\text{annealing}}+1} N_{\text{train}}(t_n),$$

where each  $N_{\text{train}}(t_n)$  satisfies the requirement (B9). The annealing times  $\{t_n\}_{n=1}^{N_{\text{annealing}}+1}$  are defined such that  $t_1 \equiv 0$  and  $t_{n+1} \equiv t_n + \delta t_n$ . Here,  $\delta t_n$  satisfies

$$\delta t_n = \mathcal{O}\left(\frac{\epsilon g(t_n)}{\text{poly}(N)}\right). \quad (\text{B10})$$

We also consider  $N_{\text{annealing}}$  the smallest integer such that  $t_{N_{\text{annealing}}} + \delta t_{N_{\text{annealing}}} \geq 1$ , in this case, we define  $t_{N_{\text{annealing}}+1} \equiv 1$ , indicating the end of annealing. Thus,  $N_{\text{annealing}}$  is the total number of annealing steps. Taking this definition into account, then one can show that

$$N_{\text{annealing}} \leq \frac{1}{\min_{\{t_n\}}(\delta t_n)} + 1.$$

Using Eqs. (B7) and (B9) and the previous inequality,  $N_{\text{steps}}$  can be bounded from above as:

$$\begin{aligned} N_{\text{steps}} &\leq (N_{\text{annealing}} + 1) \max_{\{t_n\}}(N_{\text{train}}(t_n)) \\ &\leq \left(\frac{1}{\min_{\{t_n\}}(\delta t_n)} + 2\right) \max_{\{t_n\}}(N_{\text{train}}(t_n)) \\ &\leq \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon^2 \min_{\{t_n\}}(g(t_n))^2}\right), \end{aligned}$$

where the transition from line 2 to line 3 is valid for a sufficiently small  $\epsilon$  and  $\min_{\{t_n\}}(g(t_n))$ . Furthermore,  $N_{\text{steps}}$  can also be bounded from below as:

$$N_{\text{steps}} \geq \max_{\{t_n\}}(N_{\text{train}}(t_n)) = \mathcal{O}\left(\frac{\text{poly}(N)}{\epsilon \min_{\{t_n\}}(g(t_n))}\right). \quad (\text{B11})$$

Note that the minimum in the previous two bounds are taken over all the annealing times  $t_n$  where  $1 \leq n \leq N_{\text{annealing}} + 1$ .

In this derivation of the bound on  $N_{\text{steps}}$ , we have assumed that the ground state of  $\hat{H}_{\text{target}}$  is non-degenerate, so that the gap does not vanish at the end of annealing (i.e.,  $t = 1$ ). In the case of degeneracy of the target ground state, we can define the gap  $g(t)$  by considering the lowest energy level that does not lead to the degenerate ground state.

It is also worth noting that the assumptions of this derivation can be further expanded and improved. In particular, the gradients of  $E(\boldsymbol{\lambda}, t)$  are computed stochastically (see Methods Sec. VC), as opposed to our assumption (A4) where the gradients are assumed to be known exactly. To account for noisy gradients, it is possible to use convergence bounds of stochastic gradient descent [47, 58] to estimate a bound on the number of gradient descent steps. Second-order optimization methods such as stochastic reconfiguration/natural gradient [59, 60] can potentially show a significant advantage over first-order optimization methods, in terms of scaling with the minimum gap of the time-dependent Hamiltonian  $\hat{H}(t)$ .

## Appendix C: Default Hyperparameters

In this Appendix, we summarize the architectures and the hyperparameters of the simulations performed in this paper, as shown in Tab. I. The latter has shown to yield good performance, while we believe that a more advanced study of the hyperparameters can result in optimal results. We also note that in this paper, VQA and VCA were run using a single GPU workstation for each simulation, while SQA and SA were performed on a multi-core CPU.

## Appendix D: Benchmarking Recurrent neural network cells

To show the advantage of tensorized RNNs over vanilla RNNs, we benchmark these architectures on the task of finding the ground state of the uniform ferromagnetic Ising chain (i.e.,  $J_{i,i+1} = 1$ ) with  $N = 100$  spins at the critical point (i.e., no annealing is employed). Since the couplings in this model are site-independent, we choose the parameters of the model to be also site-independent. In Fig. 9(a), we plot the energy variance per site  $\sigma^2$  (see Eq. (24)) against the number of gradient descent steps. Here  $\sigma^2$  is a good indicator of the quality of the optimized wave function [59, 61, 62]. The results show that the tensorized RNN wave function can achieve both a lower estimate of the energy variance and a faster convergence.

For the disordered systems studied in this paper, we set the weights  $T_n, U_n$  and the biases  $\mathbf{b}_n, \mathbf{c}_n$  (in Eqs. (16) and (17)) to be site-dependent. To demonstrate the benefit of using site-dependent over site-independent parameters when dealing with disordered systems, we benchmark both architectures on the task of finding the ground state of the disordered Ising chain with random discrete couplings  $J_{i,i+1} = \pm 1$  at the critical point, i.e., with a transverse field  $\Gamma = 1$ . We show the results in Fig. 9(b) and find that site-dependent parameters lead to a better performance in terms of the energy variance per spin.

Furthermore, we equally show the advantage of a dilated RNN ansatz compared to a tensorized RNN ansatz. We train both of them for the task of finding the minimum of the free energy of the Sherrington-Kirkpatrick model with  $N = 20$  spins and at temperature  $T = 1$ , as explained in Methods Sec. VB. Both RNNs have a comparable number of parameters (66400 parameters for the tensorized RNN and 59240 parameters for the dilated RNN). Interestingly, in Fig. 9(c), we find that the dilated RNN supersedes the tensorized RNN with almost an order of magnitude difference in term of the free energy variance per spin defined in Eq. (21). Indeed, this result suggests that the mechanism of skip connections allows dilated RNNs to capture long-term dependencies more efficiently compared to tensorized RNNs.

Figures	Parameter	Value
Figs. 3 and 8	Architecture Number of memory units Number of samples Initial magnetic field for VQA Initial temperature for VCA Learning rate Warmup steps Number of random instances	Tensorized RNN wave function with no-weight sharing $d_h = 40$ $N_s = 50$ $\Gamma_0 = 2$ $T_0 = 1$ $\eta = 5 \times 10^{-4}$ $N_{\text{warmup}} = 1000$ $N_{\text{instances}} = 25$
Fig. 4	Architecture Number of memory units Number of samples Initial magnetic field Initial temperature Learning rate Number of warmup steps Number of random instances	2D tensorized RNN wave function with no weight-sharing $d_h = 40$ $N_s = 25$ $\Gamma_0 = 1$ (for SQA, VQA and RVQA) $T_0 = 1$ (for SA, VCA and RVQA) $\eta = 10^{-4}$ $N_{\text{warmup}} = 1000$ for $10 \times 10$ and $N_{\text{warmup}} = 2000$ for $40 \times 40$ $N_{\text{instances}} = 25$
Figs. 5(a) and (d)	Architecture Number of memory units Number of samples Initial temperature Initial magnetic field Learning rate Number of warmup steps Number of random instances	Dilated RNN wave function with no weight-sharing $d_h = 40$ $N_s = 50$ $T_0 = 2$ (for SA and VCA) $\Gamma_0 = 2$ (for SQA) $\eta = 10^{-4}$ $N_{\text{warmup}} = 2000$ $N_{\text{instances}} = 25$
Figs. 5(b), (c), (e) and (f)	Architecture Number of memory units Number of samples Initial temperature Initial magnetic field Learning rate Number of warmup steps Number of random instances	Dilated RNN wave function with no weight-sharing $d_h = 20$ $N_s = 50$ $T_0 = 1$ (for SA and VCA) $\Gamma_0 = 1$ (for SQA) $\eta = 10^{-4}$ $N_{\text{warmup}} = 1000$ $N_{\text{instances}} = 25$
Fig. 7	Architecture Number of memory units Number of samples Initial temperature Initial magnetic field Learning rate Number of warmup steps	Tensorized RNN wave function with weight sharing $d_h = 20$ $N_s = 50$ $T_0 = 2$ $\Gamma_0 = 2$ $\eta = 10^{-3}$ $N_{\text{warmup}} = 1000$
Figs. 9(a) and (b)	Architecture Number of memory units Number of samples Learning rate	RNN wave function $d_h = 50$ $N_s = 50$ $\eta = 10^{-3}$ for Fig. 9(a) and $\eta = 5 \times 10^{-4}$ for Fig. 9(b)
Fig. 9(c)	Architecture Number of memory units of dilated RNN Number of memory units of tensorized RNN Number of samples Learning rate	RNN wave function with no-weight sharing $d_h = 20$ $d_h = 40$ $N_s = 100$ $\eta = 10^{-4}$

Table I. Hyperparameters used to obtain the results reported in this paper. Note that the number of samples stands for the batch size used to train the RNN.

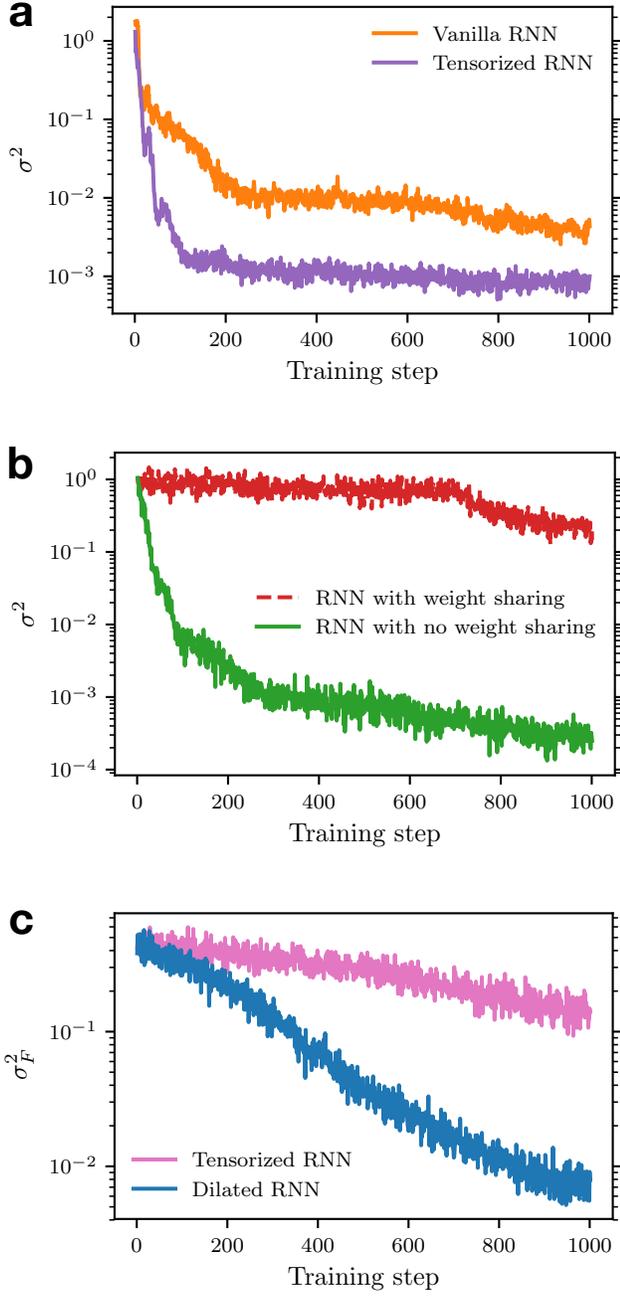


Figure 9. Energy (or Free energy) variance per spin  $\sigma^2$  vs the number of training steps. (a) We compare tensorized and vanilla RNN ansatzes both with weight sharing across sites on the uniform ferromagnetic Ising chain at the critical point with  $N = 100$  spins. (b) Comparison between a tensorized RNN with and without weight sharing, trained to find the ground state of the random Ising chain with discrete disorder ( $J_{i,i+1} = \pm 1$ ) at criticality with  $N = 20$  spins. (c) Comparison between a tensorized RNN and dilated RNN ansatzes, both with no weight sharing, trained to find the Sherrington-Kirkpatrick model's equilibrium distribution with  $N = 20$  spins at temperature  $T = 1$ .

- [1] Andrew Lucas, “Ising formulations of many np problems,” *Front. Phys.* **2**, 5 (2014).
- [2] F Barahona, “On the computational complexity of ising spin glass models,” *Journal of Physics A: Mathematical and General* **15**, 3241–3253 (1982).
- [3] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, “Optimization by simulated annealing,” *Science* **220**, 671–680 (1983).
- [4] C Koulamas, SR Antony, and R Jaen, “A survey of simulated annealing applications to operations research problems,” *Omega* **22**, 41 – 56 (1994).
- [5] Bruce Hajek, “A tutorial survey of theory and applications of simulated annealing,” in *1985 24th IEEE Conference on Decision and Control* (1985) pp. 755–760.
- [6] D.I. Svergun, “Restoring low resolution structure of biological macromolecules from solution scattering using simulated annealing,” *Biophysical Journal* **76**, 2879 – 2886 (1999).
- [7] David S. Johnson, Cecilia R. Aragon, Lyle A. McGeoch, and Catherine Schevon, “Optimization by simulated annealing: An experimental evaluation; part ii, graph coloring and number partitioning,” *Operations Research* **39**, 378–406 (1991).
- [8] M. A. Abido, “Robust design of multimachine power system stabilizers using simulated annealing,” *IEEE Transactions on Energy Conversion* **15**, 297–304 (2000).
- [9] Torsten Karzig, Armin Rahmani, Felix von Oppen, and Gil Refael, “Optimal control of majorana zero modes,” *Phys. Rev. B* **91**, 201404 (2015).
- [10] Georges Gielen, Herman Walscharts, and Willy Sansen, “Analog circuit design optimization based on symbolic simulation and simulated annealing,” in *ESSCIRC ’89: Proceedings of the 15th European Solid-State Circuits Conference* (1989) pp. 252–255.
- [11] Giuseppe E. Santoro, Roman Martoňák, Erio Tosatti, and Roberto Car, “Theory of quantum annealing of an ising spin glass,” *Science* **295**, 2427–2430 (2002).
- [12] J. Brooke, D. Bitko, T. F. Rosenbaum, and G. Aeppli, “Quantum annealing of a disordered magnet,” *Science* **284**, 779–781 (1999).
- [13] Debasis Mitra, Fabio Romeo, and Alberto Sangiovanni-Vincentelli, “Convergence and finite-time behavior of simulated annealing,” *Advances in Applied Probability* **18**, 747–771 (1986).
- [14] Daniel Delahaye, Supatcha Chaimatanan, and Marcel Mongeau, “Simulated annealing: From basics to applications,” in *Handbook of Metaheuristics*, edited by Michel Gendreau and Jean-Yves Potvin (Springer International Publishing, Cham, 2019) pp. 1–35.
- [15] Ilya Sutskever, James Martens, and Geoffrey Hinton, “Generating text with recurrent neural networks,” in *Proceedings of the 28th International Conference on International Conference on Machine Learning, ICML’11* (Omnipress, Madison, WI, USA, 2011) p. 1017–1024.
- [16] Hugo Larochelle and Iain Murray, “The neural autoregressive distribution estimator,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, Vol. 15, edited by Geoffrey Gordon, David Dunson, and Miroslav Dudík (JMLR Workshop and Conference Proceedings, Fort Lauderdale, FL, USA, 2011) pp. 29–37.
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” (2017), arXiv:1706.03762 [cs.CL].
- [18] Dian Wu, Lei Wang, and Pan Zhang, “Solving statistical mechanics using variational autoregressive networks,” *Physical Review Letters* **122** (2019), 10.1103/physrevlett.122.080602.
- [19] Or Sharir, Yoav Levine, Noam Wies, Giuseppe Carleo, and Amnon Shashua, “Deep autoregressive models for the efficient variational simulation of many-body quantum systems,” *Physical Review Letters* **124** (2020), 10.1103/physrevlett.124.020503.
- [20] Mohamed Hibat-Allah, Martin Ganahl, Lauren E. Hayward, Roger G. Melko, and Juan Carrasquilla, “Recurrent neural network wave functions,” *Physical Review Research* **2** (2020), 10.1103/physrevresearch.2.023358.
- [21] Christopher Roth, “Iterative retraining of quantum spin models using recurrent neural networks,” (2020), arXiv:2003.06228 [physics.comp-ph].
- [22] R.P. Feynman, *Statistical Mechanics: A Set of Lectures*, Advanced Books Classics (Avalon Publishing, 1998).
- [23] Philip M. Long and Rocco A. Servedio, “Restricted boltzmann machines are hard to approximately evaluate or simulate,” in *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML’10* (Omnipress, Madison, WI, USA, 2010) p. 703–710.
- [24] Sergio Boixo, Troels F Rønnow, Sergei V Isakov, Zhihui Wang, David Wecker, Daniel A Lidar, John M Martinis, and Matthias Troyer, “Evidence for quantum annealing with more than one hundred qubits,” *Nat. Phys.* **10**, 218–224 (2014).
- [25] Tadashi Kadowaki and Hidetoshi Nishimori, “Quantum annealing in the transverse ising model,” *Physical Review E* **58**, 5355–5363 (1998).
- [26] M. Born and V. Fock, “Beweis des adiabatenatzes,” *Zeitschrift für Physik* **51**, 165–180 (1928).
- [27] Glen Bigan Mbeng, Lorenzo Privitera, Luca Arceci, and Giuseppe E. Santoro, “Dynamics of simulated quantum annealing in random ising chains,” *Phys. Rev. B* **99**, 064201 (2019).
- [28] Nilan Norris, “The standard errors of the geometric and harmonic means and their application to index numbers,” *The Annals of Mathematical Statistics* **11**, 445–448 (1940).
- [29] Tommaso Zanca and Giuseppe E. Santoro, “Quantum annealing speedup over simulated annealing on random ising chains,” *Phys. Rev. B* **93**, 224431 (2016).
- [30] “<https://software.cs.uni-koeln.de/spinglass/>,” .
- [31] Neil G Dickson, MW Johnson, MH Amin, R Harris, F Altomare, AJ Berkley, P Bunyk, J Cai, EM Chapple, P Chavez, *et al.*, “Thermally assisted quantum annealing of a 16-qubit problem,” *Nature communications* **4**, 1–6 (2013).
- [32] Joseph Gomes, Keri A. McKiernan, Peter Eastman, and Vijay S. Pande, “Classical quantum optimization with neural network quantum states,” (2019), arXiv:1910.10675 [cond-mat.dis-nn].
- [33] Semyon Sinchenko and Dmitry Bazhanov, “The deep

- learning and statistical physics applications to the problems of combinatorial optimization,” (2019), arXiv:1911.10680 [cond-mat.dis-nn].
- [34] Tianchen Zhao, Giuseppe Carleo, James Stokes, and Shравan Veerapaneni, “Natural evolution strategies and quantum approximate optimization,” (2020), arXiv:2005.04447 [quant-ph].
- [35] Roman Martoňák, Giuseppe E. Santoro, and Erio Tosatti, “Quantum annealing by the path-integral monte carlo method: The two-dimensional random ising model,” *Phys. Rev. B* **66**, 094203 (2002).
- [36] M Mezard, G Parisi, and M Virasoro, *Spin Glass Theory and Beyond* (WORLD SCIENTIFIC, 1986) <https://www.worldscientific.com/doi/pdf/10.1142/0271>.
- [37] David Sherrington and Scott Kirkpatrick, “Solvable model of a spin-glass,” *Phys. Rev. Lett.* **35**, 1792–1796 (1975).
- [38] Firas Hamze, Jack Raymond, Christopher A. Pattison, Katja Biswas, and Helmut G. Katzgraber, “Wishart planted ensemble: A tunably rugged pairwise ising model with a first-order phase transition,” *Physical Review E* **101** (2020), 10.1103/physreve.101.052102.
- [39] Kyle Mills, Pooya Ronagh, and Isaac Tamblyn, “Controlled online optimization learning (cool): Finding the ground state of spin hamiltonians with reinforcement learning,” (2020), arXiv:2003.00011 [physics.comp-ph].
- [40] Yoshua Bengio, Andrea Lodi, and Antoine Prouvost, “Machine learning for combinatorial optimization: A methodological tour d’horizon,” *European Journal of Operational Research* (2020), <https://doi.org/10.1016/j.ejor.2020.07.063>.
- [41] Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning* (MIT Press, 2016) <http://www.deeplearningbook.org>.
- [42] Richard Kelley, “Sequence modeling with recurrent tensor networks,” (2016).
- [43] Shiyu Chang, Yang Zhang, Wei Han, Mo Yu, Xiaoxiao Guo, Wei Tan, Xiaodong Cui, Michael Witbrock, Mark Hasegawa-Johnson, and Thomas S. Huang, “Dilated recurrent neural networks,” (2017), arXiv:1710.02224 [cs.AI].
- [44] Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE Transactions on Neural Networks* **5**, 157–166 (1994).
- [45] Salah El Hihi and Yoshua Bengio, “Hierarchical recurrent neural networks for long-term dependencies,” in *Advances in Neural Information Processing Systems 8*, edited by D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo (MIT Press, 1996) pp. 493–499.
- [46] G. Vidal, “Class of quantum many-body states that can be efficiently simulated,” *Physical Review Letters* **101** (2008), 10.1103/physrevlett.101.110501.
- [47] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” (2014), arXiv:1412.6980 [cs.LG].
- [48] Sergey Bravyi, David P. Divincenzo, Roberto Oliveira, and Barbara M. Terhal, “The complexity of stoquastic local hamiltonian problems,” *Quantum Info. Comput.* **8**, 361–385 (2008).
- [49] I. Ozfidan, C. Deng, A.Y. Smirnov, T. Lanting, R. Harris, L. Swenson, J. Whittaker, F. Altomare, M. Babcock, C. Baron, A.J. Berkley, K. Boothby, H. Christiani, P. Bunyk, C. Enderud, B. Evert, M. Hager, A. Haja, J. Hilton, S. Huang, E. Hoskinson, M.W. Johnson, K. Jooya, E. Ladizinsky, N. Ladizinsky, R. Li, A. MacDonald, D. Marsden, G. Marsden, T. Medina, R. Molavi, R. Neufeld, M. Nissen, M. Norouzpour, T. Oh, I. Pavlov, I. Perminov, G. Poulin-Lamarre, M. Reis, T. Prescott, C. Rich, Y. Sato, G. Sterling, N. Tsai, M. Volkmann, W. Wilkinson, J. Yao, and M.H. Amin, “Demonstration of a nonstoquastic hamiltonian in coupled superconducting flux qubits,” *Phys. Rev. Applied* **13**, 034037 (2020).
- [50] Mohamed Hibat-Allah, Estelle M. Inack, Roger G. Melko, and Juan Carrasquilla, (Manuscript in preparation).
- [51] Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih, “Monte carlo gradient estimation in machine learning,” (2019), arXiv:1906.10652 [stat.ML].
- [52] Shi-Xin Zhang, Zhou-Quan Wan, and Hong Yao, “Automatic differentiable monte carlo: Theory and application,” (2019), arXiv:1911.09117 [physics.comp-ph].
- [53] Sei Suzuki, “Cooling dynamics of pure and random ising chains,” *Journal of Statistical Mechanics: Theory and Experiment* **2009**, P03032 (2009).
- [54] Jacek Dziarmaga, “Dynamics of a quantum phase transition in the random ising model: Logarithmic dependence of the defect density on the transition rate,” *Phys. Rev. B* **74**, 064416 (2006).
- [55] Tommaso Caneva, Rosario Fazio, and Giuseppe E. Santoro, “Adiabatic quantum dynamics of a random ising chain across its quantum critical point,” *Phys. Rev. B* **76**, 144427 (2007).
- [56] Sandro Sorella and Federico Becca, *SISSA Lecture notes on Numerical methods for strongly correlated electrons (Sec. 1.3)* (2016).
- [57] Yurii Nesterov, “Smooth convex optimization,” in *Lectures on Convex Optimization* (Springer International Publishing, Cham, 2018) pp. 59–137.
- [58] Mark Schmidt, Nicolas Le Roux, and Francis Bach, “Minimizing finite sums with the stochastic average gradient,” (2013), arXiv:1309.2388 [math.OC].
- [59] F. Becca and S. Sorella, *Quantum Monte Carlo Approaches for Correlated Systems* (Cambridge University Press, 2017).
- [60] Shun-ichi Amari, “Natural gradient works efficiently in learning,” *Neural Computation* **10**, 251–276 (1998), <https://doi.org/10.1162/089976698300017746>.
- [61] Claudius Gros, “Criterion for a good variational wave function,” *Phys. Rev. B* **42**, 6835–6838 (1990).
- [62] Roland Assaraf and Michel Caffarel, “Zero-variance zero-bias principle for observables in quantum monte carlo: Application to forces,” *The Journal of Chemical Physics* **119**, 10536–10552 (2003).