# Probabilistic Computers for Neural Quantum States

Shuvro Chowdhury,[1] Jasper Pieterse,[2] Navid Anjum Aadit,[1] Johan H. Mentink,[2] and Kerem Y. Camsari[1]

[1]*Department of Electrical and Computer Engineering,*
*University of California, Santa Barbara, Santa Barbara, CA 93106, USA*
[2]*Radboud University, Institute for Molecules and Materials, Heyendaalseweg 135, Nijmegen, The Netherlands*
(Dated: January 1, 2026)

Neural quantum states efficiently represent many-body wavefunctions with neural networks, but the cost of Monte Carlo sampling limits their scaling to large system sizes. Here we address this challenge by combining sparse Boltzmann machine architectures with probabilistic computing hardware. We implement a probabilistic computer on field-programmable gate arrays (FPGAs) and use it as a fast sampler for energy-based neural quantum states. For the two-dimensional transverse-field Ising model at criticality, we obtain accurate ground-state energies for lattices up to $80{\times}80$ (6400 spins) using a custom multi-FPGA cluster. Furthermore, we introduce a dual-sampling algorithm to train deep Boltzmann machines, replacing intractable marginalization with conditional sampling over auxiliary layers. This enables the training of sparse deep models and improves parameter efficiency relative to shallow networks. Using this algorithm, we train deep Boltzmann machines for a system with $35 \times 35$ (1225 spins). Together, these results demonstrate that probabilistic hardware can overcome the sampling bottleneck in variational simulation of quantum many-body systems, opening a path to larger system sizes and deeper variational architectures.

## I. Introduction

Many-body quantum simulation drives progress in condensed matter physics, quantum chemistry, and quantum information, yet accurate classical simulation remains challenging as system size grows. Established approaches such as quantum Monte Carlo (QMC) and tensor-network (TN) methods have achieved high precision in important regimes, but each faces fundamental limitations: sign problems for generic QMC and unfavorable entanglement scaling for tensor networks in two dimensions and near criticality [1–3].

Neural Quantum States (NQS), introduced by Carleo and Troyer [4], offer a different approach (Fig. 1a,b). By parameterizing many-body wavefunctions with neural networks, NQS can represent complex correlations in a compact and systematically improvable form [5, 6]. This framework has enabled progress in two-dimensional quantum systems beyond the reach of exact diagonalization [7, 8].

Since their inception, NQS architectures have diversified. While early work relied on Restricted Boltzmann Machines (RBMs), subsequent studies adopted convolutional networks to exploit translational symmetry [9], recurrent neural networks [10], transformer-based models [11–13] and more recently foundation models [14], which significantly increased computational costs compared to RBMs. Despite this diversity, RBMs remain a canonical ansatz and continue to be used in studies of topological phases and fracton models [15, 16].

Early NQS applications demonstrated simulations of systems with hundreds of spins, far beyond what is accessible with exact diagonalization. Scaling to larger system sizes remains an active area of research, with advances in software frameworks [17–20], algorithms [21–23], and high-performance computing [24–26], now reaching systems on the order of $10^3$ spins.

However, increasing system size by another order of magnitude remains challenging. For variational Monte Carlo (VMC), a widely used and expressive approach [27], scaling is constrained by the cost of Markov Chain Monte Carlo (MCMC) sampling required to estimate observables and stochastic parameter gradients. This bottleneck persists even for relatively simple architectures such as RBMs and becomes more severe for deeper or more structured models.

Probabilistic hardware provides a promising route to overcome this limitation. The joint energy of an RBM corresponds to an Ising Hamiltonian on a bipartite graph, allowing nonlocal correlations to be mediated through local pairwise interactions that map naturally onto probabilistic bits (p-bits) [28, 29]. P-bits are classical stochastic units that fluctuate between two logic states and can be implemented with massive parallelism and low-precision arithmetic [30]. Recent theoretical work predicts large sampling speedups for NQS when using probabilistic hardware [31], and experimental demonstrations have shown orders-of-magnitude acceleration for frustrated Ising models using hardware-accelerated convolutional RBMs [32]. Prior p-bit architectures have established that local connectivity enables sparse Boltzmann machines to map efficiently onto hardware while avoiding communication bottlenecks [33, 34].

Here, we implement a probabilistic computer (p-computer) using field-programmable gate arrays (FPGAs) to accelerate sampling for neural quantum states (Fig. 1c). Rather than emulating stochasticity in software, we map sparse Boltzmann machines directly onto the FPGA fabric, exploiting spatial parallelism and low-precision arithmetic. Enforcing local connectivity enables efficient on-chip realization and avoids the routing congestion inherent to dense networks.

This work makes two complementary contributions. First, we map Further Restricted Boltzmann Machine (FRBM) − a sparse variant of RBM − onto probabilistic hardware. Using hardware-parallel sampling on a custom multi-FPGA cluster, we scale to lattices as large as $80 \times 80$ spins (6400) and obtain ground-state energies within chemical accuracy of variational benchmarks for the two-dimensional transverse-field Ising model at criticality.
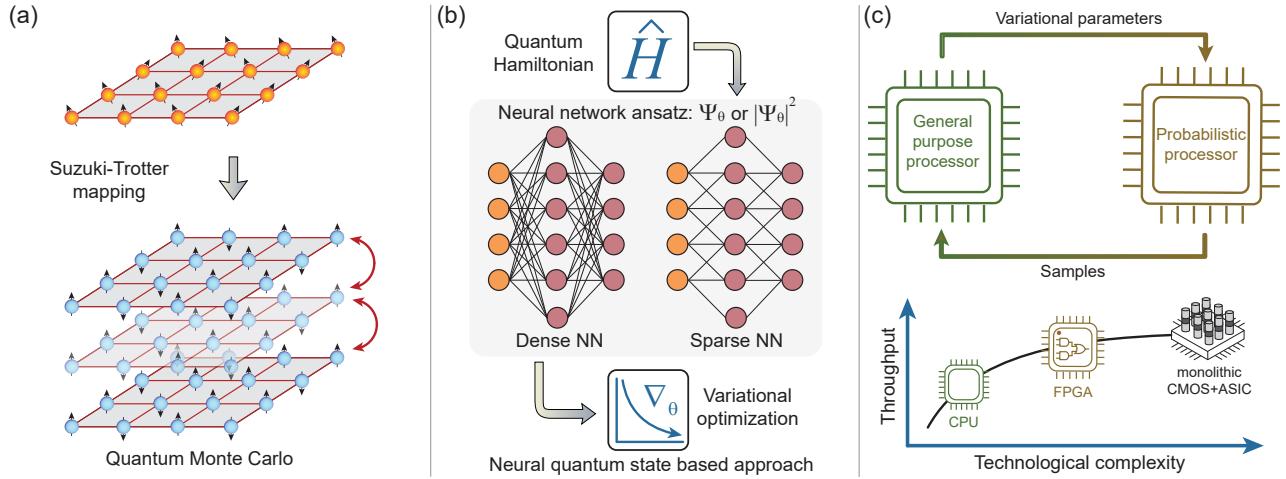
**Fig. 1**. **Probabilistic computers for accelerating variational simulations of quantum systems.** (a) Conventional approaches such as path-integral quantum Monte Carlo (QMC) map quantum systems to classical statistical models in higher dimensions via the Suzuki-Trotter decomposition, but often suffer from sign problems or unfavorable scaling. (b) Neural Quantum States (NQS) represent the many-body wavefunction with neural networks ansätze. While many NQS architectures rely on dense or structured networks, sparse locally connected models are particularly well suited for mapping onto probabilistic hardware used in this work. Ground states are obtained through variational optimization of the quantum energy, which requires repeated Monte Carlo sampling from the neural network model. (c) Sampling is offloaded to a dedicated probabilistic processor built from p-bits. A general-purpose processor updates variational parameters while a p-bit array provides high-throughput samples. Throughput increases as implementations progress from CPU emulation to FPGA-based systems and, ultimately, to dedicated CMOS hardware.

Second, we introduce a dual-sampling algorithm for deep Boltzmann machines (DBMs) that replaces intractable marginalization over auxiliary variables with conditional sampling. This reformulation enables stable training of sparse deep architectures under strict locality constraints and improves parameter efficiency relative to shallow networks. We demonstrate dual sampling for DBMs up to $35 \times 35$ spins (1225) on a single GPU.

Together, these advances address two barriers to scaling neural quantum states: probabilistic hardware alleviates the sampling bottleneck, while dual sampling makes sparse deep Boltzmann machines trainable for variational Monte Carlo.

## II. Model and Ansatz

We target the two-dimensional transverse-field Ising model (TFIM) on a square lattice with periodic boundary conditions, defined by the Hamiltonian [37]:

$$\hat{H} = -J \sum_{\langle i,j \rangle} \hat{\sigma}_i^z \hat{\sigma}_j^z - \Gamma_x \sum_i \hat{\sigma}_i^x \quad (1)$$

where $\langle i, j \rangle$ denotes unique nearest-neighbor bonds, $J$ is the coupling strength, and $\Gamma_x$ is the transverse field. We define the many-body basis states $S = \{s_1, s_2, \ldots, s_N\}$ where $s_i \in \{+1, -1\}$ and parameterize the squared wavefunction amplitude $|\Psi_\theta(S)|^2$ through the Boltzmann distribution of a probabilistic network [4]:

$$|\Psi_\theta(S)|^2 = P_\theta(S) = \frac{1}{Z_\theta} \sum_h e^{-E_\theta(S,h)} \quad (2)$$

where $h = \{h_1, \ldots, h_M\}$ with $h_j \in \{+1, -1\}$. This mapping relies on the Perron-Frobenius theorem, which

guarantees that the ground state of stoquastic Hamiltonians like the TFIM can be chosen to have non-negative real amplitudes [38]. The quantum probability distribution $|\Psi|^2$ then coincides with the classical distribution generated by p-bits. As a result, the p-computer's stochastic fluctuations directly generate samples from $|\Psi_\theta(S)|^2$ without computing amplitudes explicitly. Extending this approach to non-stoquastic systems with nontrivial sign structures would require complex parameters or phase networks [39–41].

The variational ground-state energy of a neural quantum state (NQS) $\Psi_\theta(S)$ is evaluated as

$$E = \frac{\langle \Psi_\theta | H | \Psi_\theta \rangle}{\langle \Psi_\theta | \Psi_\theta \rangle} = \mathbb{E}_{S \sim |\Psi_\theta|^2}[E_{\text{loc}}(S)] \quad (3)$$

where the local energy is defined as

$$E_{\text{loc}}(S) = \sum_{S'} \frac{\Psi_\theta(S')}{\Psi_\theta(S)} \langle S | H | S' \rangle \quad (4)$$

Training the neural quantum state corresponds to minimizing this variational energy $E(\theta)$ with respect to the network parameters $\theta$. Expectation values are obtained by Monte Carlo sampling of configurations $S$ drawn from the probability distribution $|\Psi_\theta(S)|^2$. For the transverse-field Ising model, the off-diagonal term $-\Gamma_x \sum_i \sigma_i^x$ couples a configuration $S$ to configurations $S^{(i)}$ that differ by a single spin flip at site $i$. As a result, evaluating $E_{\text{loc}}(S)$ requires computing wavefunction ratios of the form $\Psi_\theta(S^{(i)})/\Psi_\theta(S)$ for sampled configurations. These ratios also appear in estimators of physical observables such as $\langle \sigma_i^x \rangle$. For shallow architectures such as RBMs, these wavefunction ratios can be evaluated analytically. However, for deep Boltzmann machines, the wavefunction amplitude involves an intractable
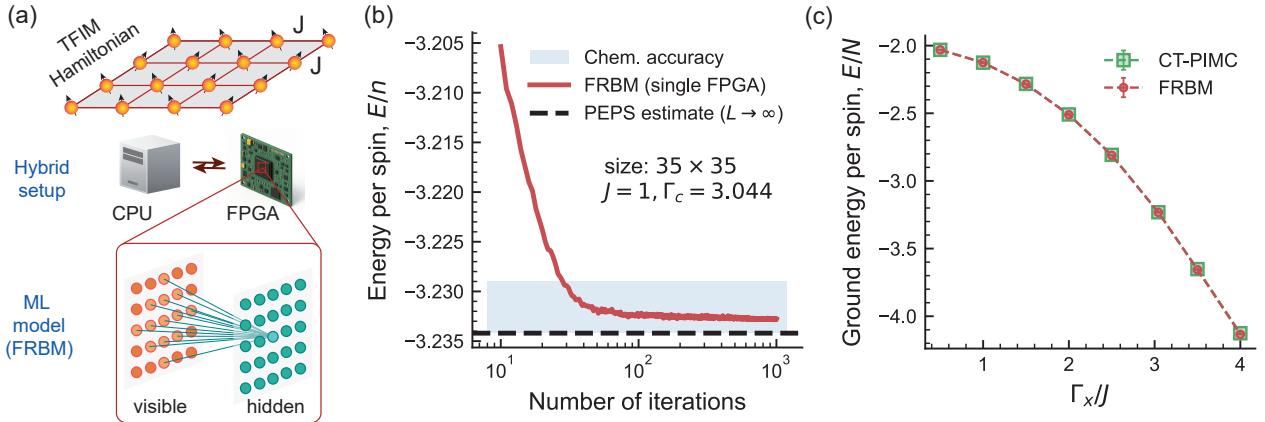
**Fig. 2**. **Single-FPGA results for the 2D transverse-field Ising model.** (a) Problem and hybrid setup: a 2D square lattice with periodic boundaries, simulated using a CPU-FPGA platform where the FPGA samples from a Further Restricted Boltzmann Machine (FRBM) with local connectivity ($k = 2$, corresponding to 13 neighbors). (b) Training convergence for a $35 \times 35$ lattice (1225 spins) at the critical field $\Gamma_c/J = 3.044$. The energy per spin reaches chemical accuracy (blue shaded band: relative error $|\Delta E/E_{\mathrm{ref}}| \leq 1.6 \times 10^{-3}$, $E_{\mathrm{ref}}$ is a variational Projected Entangled Pair States (PEPS) benchmark estimate at the thermodynamic limit ($L \to \infty$) from Ref. [35]) within $\approx$100 iterations. (c) Ground-state energy per spin versus transverse field, interpolating between the ferromagnetic limit ($E/N \to -2$ as $\Gamma_x \to 0$) and the field polarised limit ($E/N \to -\Gamma_x$ as $\Gamma_x \to \infty$). For validation, we also compare against Continuous-Time Path Integral Monte Carlo simulation ($\beta = 64$ and $10^5$ Monte Carlo sweeps) using the code provided in [36]. FRBM results are obtained from an average over $10^6$ samples of a single run after the training is completed. Error bars represent standard error via blocking (50 bins) and are smaller than the symbol size.

marginalization over hidden units, rendering direct evaluation of such ratios prohibitively expensive. To overcome this bottleneck, we introduce in Section VI a dual sampling strategy that estimates wavefunction ratios via conditional sampling over hidden variables while preserving unbiased estimators of the local energy.

### III. Sparse Architecture and Hardware Implementation

The physical implementation (Fig. 2a) uses a hybrid probabilistic-classical architecture [34, 42]: a Xilinx Alveo U250 FPGA serves as the probabilistic processor while a host CPU handles parameter optimization. The U250 scales to large quantum systems through the sparse architecture described below.

Central to this approach is the Further Restricted Boltzmann Machine (FRBM) [43] with strictly local connectivity. Standard all-to-all RBMs require $\mathcal{O}(N^2)$ connections, creating routing congestion on hardware. The FRBM instead defines the network energy $E_\theta(S, h)$ as:

$$E_\theta(S, h) = -\sum_i a_i s_i - \sum_j b_j h_j - \sum_{\langle i,j \rangle_k} W_{ij} s_i h_j \quad (5)$$

where $h = \{h_1, h_2, \ldots, h_M\}$ represents the hidden spins, $a, b, W$ are the variational parameters, and $\langle i, j \rangle_k$ restricts connections to nodes within Euclidean distance $k$. This reduces wiring complexity to $\mathcal{O}(N)$ for fixed $k$. Throughout this work we use $k = 2$, corresponding to 13 neighbors per spin. This sparsity maps directly onto the FPGA while preserving analytical tractability for wavefunction amplitudes, local energies, and gradients [4, 44].

To maximize p-bit density on the Alveo U250, we use 10-bit fixed-point arithmetic (1 sign, 6 integer, 3 fractional

bits) for weights and states, reducing logic utilization per neuron and enabling greater parallelism than floating-point implementations. For uncorrelated stochasticity, we integrate the xoshiro pseudo-random number generator [45], which provides sufficient statistical quality to avoid sampling bias while maintaining a compact footprint [46]. These choices allow a large ensemble of p-bits on a single device without compromising the precision needed for accurate energy estimation. On the host CPU, the evaluation of local energies, gradients, and stochastic reconfiguration updates is performed in single-precision floating point (FP32), providing sufficient numerical accuracy while maintaining high throughput.

Parameter optimization uses the Stochastic Reconfiguration (SR) algorithm [4, 47]. The variational optimization loop is split between the FPGA and the host CPU: given a set of parameters, the FPGA generates spin configurations which are transferred to the host CPU for parameter updates, and the process is repeated until the variational energy $E(\theta)$ in Eq. (3) converges. To avoid inverting the dense Fisher Information Matrix, we use an iterative Conjugate Gradient (CG) solver with matrix-free matrix-vector products [48], which scales efficiently with network size. Because the CG steps involve simple local operations, future implementations could integrate parameter optimization directly onto the FPGA.

### IV. Validation on single-FPGA

We validated the method on a $35 \times 35$ lattice (1225 spins, 2450 p-bits). As shown in Fig. 2b, the variational energy per spin converges to the benchmark ground-state value [35], reaching chemical accuracy (relative error $|\Delta E/E_{\mathrm{ref}}| \leq 1.6 \times 10^{-3}$) within approximately 100
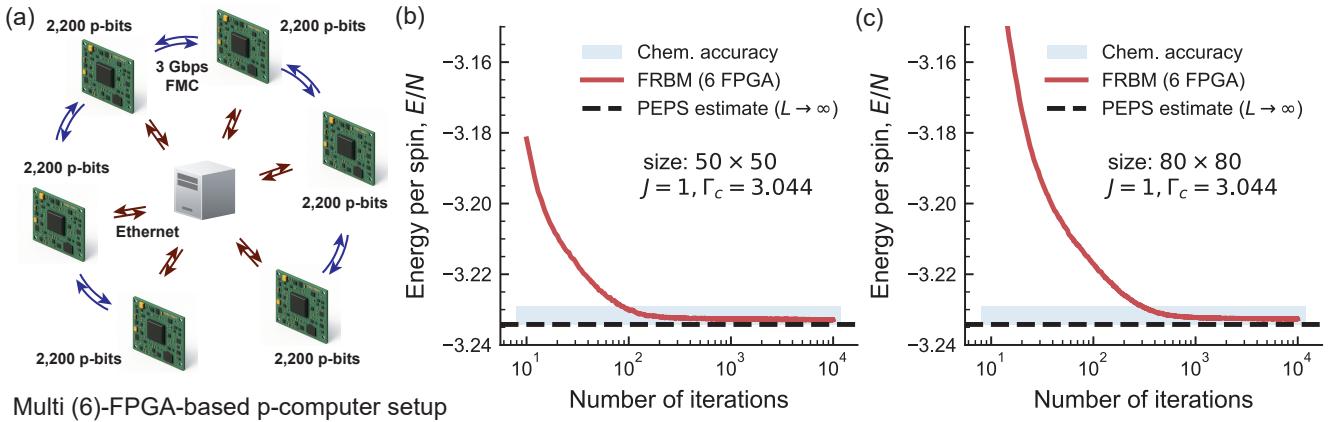
**Fig. 3**. **Multi-FPGA results for large-scale neural quantum states.** (a) Six-FPGA probabilistic computer with boards connected via 3 Gbps FMC links and coordinated over Ethernet; each FPGA hosts up to 2200 p-bits. (b) Training convergence for a $50 \times 50$ lattice (2500 spins) at the critical field $\Gamma_c/J = 3.044$. (c) Same for an $80 \times 80$ lattice (6400 spins), showing empirically that convergence within chemical accuracy (blue shaded region) is maintained as system size increases. PEPS benchmark estimate at the thermodynamic limit is taken from Ref. [35].

optimization iterations. Each iteration involves $10^4$ MCMC samples (see Supplementary Section E).

Profiling the training loop (Supplementary Section A) shows that sampling consumes less than 5% of total wall-clock time on the FPGA, despite executing $3 \times 10^6$ Monte Carlo sweeps (MCS) per iteration. In contrast, a CPU baseline allocates 20–30% of runtime to sampling while being limited to $10^4$ MCS per iteration as same sampling budget as FPGA would make training prohibitively slow. Profiling the training loop (Supplementary Section A) shows that sampling consumes less than 5% of total wall-clock time on the FPGA, despite executing $3 \times 10^6$ Monte Carlo sweeps per iteration. In contrast, a CPU baseline limited to $10^4$ sweeps per iteration already spends 20–30% of runtime on sampling; matching the FPGA's throughput would be prohibitively slow. This reduction in sampling time aligns with theoretical predictions of large sampling speedups for stochastic Ising machines applied to NQS [31]. Crucially, the sampling time per sweep remains constant as system size increases, provided the FRBM fits within the resources of the probabilistic chip [42, 49]. With sampling latency minimized, local energy computation and CG-based gradient accumulation dominate the remaining execution time. The sparse connectivity of our architecture opens a path to offloading local-energy calculations to hardware, an optimization we leave for future work. To further validate the solution, we swept the transverse field $\Gamma_x/J$ across the critical point [50]. Fig. 2c shows the ground-state energy per spin interpolating smoothly between the ferromagnetic limit (low $\Gamma_x$) and the field-polarized limit (high $\Gamma_x$). These results indicate that sparse connectivity, together with low-precision sampling, provides sufficient accuracy for large lattices.

**V. Large-Scale Systems via Multi-FPGA Clusters**

A single FPGA suffices for moderately sized lattices, but larger systems require distributing the computation. We developed a cluster of six interconnected FPGA boards (Supplementary Section B), partitioning the FRBM with the graph partitioning tool METIS to minimize the number of boundary p-bits communicated between devices. Each FPGA updates its local p-bits using local fields and the most recently received boundary states, which are exchanged asynchronously and held constant between communication events.

Each FPGA hosts a subgraph containing approximately $850-2200$ p-bits depending on the system size. METIS partitioning yields cut fractions of 8.6% for $L = 50$ and 5.6% for $L = 80$, so that most interactions remain on-chip. The FPGAs are connected in a linear topology via 3 Gbps FMC links. Boundary p-bits are exchanged asynchronously over the FMC links while local p-bits update synchronously within each FPGA.

Strict global synchronization would limit the local p-bit update clock to 2.4 MHz for $L = 50$ and 1.2 MHz for $L = 80$. In practice, because boundary p-bits constitute a small fraction of the lattice, the system tolerates delayed boundary updates, allowing local p-bits to be overclocked to 15 MHz. A host computer loads weights and reads out states over Ethernet.

Using this architecture, we simulate the 2D TFIM at the critical point $\Gamma_c/J = 3.044$ for $50 \times 50$ (2500 spins) and $80 \times 80$ (6400 spins) lattices (Fig. 3b,c). These results show that p-bit-accelerated NQS reaches system sizes beyond those reported for CPU- and GPU-based implementations while maintaining accuracy within chemical precision of variational benchmarks.

The convergence behavior remains comparable across system sizes: for both lattices, the variational energy density converges smoothly to chemical accuracy (blue shaded region). The accuracy on the $80 \times 80$ lattice confirms that the distributed architecture and asynchronous communication do not introduce artifacts that degrade the solution.

## VI.  Deep Neural Quantum States via Dual Sampling

---

**Algorithm 1:** Dual Sampling Algorithm for Deep NQS

---

**Input:** Hamiltonian parameters, Learning rate $\eta$, Samples $N_s$ (outer), $N_c$ (inner)

**Output:** Optimized parameters $\theta$ (weights and biases)

1  **for** *each optimization step $t = 1$ to $N_{\text{iter}}$*
2     Generate $N_s$ visible configurations $v$ by sampling the physical layer.
3     **for** *each visible configuration $v$*
4         Clamp visible units to fixed state $v$.
5         Run $N_c$ Gibbs steps for auxiliary layers $(h, d)$.
6         Accumulate conditional expectations for wavefunction ratios and gradients.
7         Compute local energy $E_{\text{loc}}(v)$ by combining single-spin-flip ratios.
8     Compute natural gradients $\nabla\theta$ using Stochastic Reconfiguration (SR).
9     Update parameters: $\theta \leftarrow \theta - \eta(t)\nabla\theta$.
10 **return** Optimized $\theta$ and $E_{\text{G}}^{(N_{\text{iter}})}$.

---

While the shallow FRBM suffices for the TFIM, more complex phases with topological order or frustration require correlations beyond a finite local horizon. On hardware with strict locality constraints, shallow networks face a geometric ceiling that width alone cannot overcome. Deep Boltzmann Machines (DBMs) address this by stacking local layers to expand receptive fields, allowing global correlations to emerge while preserving hardware-compatible sparsity. Theory supports this: deep architectures can efficiently encode volume-law entanglement and states from polynomial-depth quantum circuits inaccessible to shallow RBMs [51]. Constructive proofs further show that DBMs can exactly represent ground states where shallow networks fail [52]. Extending the FRBM energy in Eq. (5), we define a sparse DBM energy by coupling the hidden layer $h$ to an additional deep layer $d$:

$$E_\theta(S, h, d) = -\sum_i a_i s_i - \sum_j b_j h_j - \sum_{\langle i,j\rangle_{k_1}} W_{ij}^{(vh)} s_i h_j$$
$$- \sum_l c_l d_l - \sum_{\langle j,l\rangle_{k_2}} h_j W_{jl}^{(hd)} d_l \qquad (6)$$

where $W^{(hd)}$ connects local neighbors in layers $h$ and $d$.

This depth precludes the analytic marginalization that makes local energy and gradient calculations tractable in RBMs. Evaluating local energy and gradients in DBMs requires nested sampling over auxiliary variables [53], which has limited the adoption of DBMs for neural quantum states [54].

We address this challenge by developing a dual sampling algorithm (Algorithm 1) that reformulates variational Monte Carlo for deep neural quantum states in terms of conditional expectations. This strategy decouples sampling of physical spins from conditional sampling of auxiliary layers, while preserving an asymptotically unbiased estimator of the variational energy.

Conceptually, this mirrors contrastive-learning-style training in energy-based models, where clamped (data-conditioned) correlations are estimated by sampling hidden variables conditioned on fixed visible data [55, 56]. In the present setting, the "data" corresponds to a fixed visible spin configuration $S$, and auxiliary variables are sampled conditionally for each such configuration. This algorithm exploits the sparse structure of DBMs and aligns naturally with fast probabilistic samplers.

Within this framework, the wavefunction amplitude ratios are expressed as conditional expectations over the auxiliary variables given a fixed visible configuration:

$$\frac{\Psi(S^{(i)})}{\Psi(S)} = \sqrt{\frac{P(S^{(i)})}{P(S)}} = \sqrt{\mathbb{E}_{(h,d)\sim P(h,d|S)}\left[e^{-\Delta E_i(S,h,d)}\right]} \qquad (7)$$

where $\Delta E_i = E(S^{(i)}, h, d) - E(S, h, d)$ is the change in the classical energy of the DBM upon flipping a visible spin in $S$ at site $i$, and the expectation is estimated via samples drawn from the conditional distribution $P(h, d|S)$.

Conditioning on the visible configuration substantially reduces variance by restricting sampling to the auxiliary variables. Reusing each of the $N_c$ conditional samples to evaluate all single-spin-flip ratios further reduces computational cost by avoiding resampling for each spin flip. This construction yields an asymptotically unbiased estimator of the wavefunction ratio that converges to the exact value in the infinite-sample limit (Supplementary Section D). The nonlinear square-root mapping from probabilities to amplitudes introduces a residual bias at finite $N_c$. To mitigate this, we apply a second-order Taylor correction to the spin-flip probability estimates, as detailed in Supplementary Algorithm S1.

We validate the approach on a $10 \times 10$ lattice, using $10^3$ inner-loop samples to estimate local energy and parameter derivatives. At this size, optimization results can be directly compared against reliable benchmarks. As shown in Fig. 4a, the variational energy converges within chemical accuracy (blue shaded region) relative to a recent PEPS benchmark [35]. For the system studied here, this empirical convergence indicates that $10^3$ conditional samples in the inner-loop suffice to maintain the correct balance between diagonal and off-diagonal contributions to the local energy, keeping stochastic reconfiguration updates well-conditioned. Supplementary Section G provides a computational complexity analysis of the algorithm.

Beyond convergence and accuracy, Fig. 4b compares representational efficiency between shallow and deep architectures. We swept the number of variational parameters $N_p$ for both the single-layer sparse RBM (blue circles) and the sparse DBM (green squares). While both architectures reach chemical accuracy, the sparse DBM achieves lower variational energies with less than half the parameters ($N_p \approx 1300$ versus $N_p \approx 3100$), consistent with theoretical results showing that
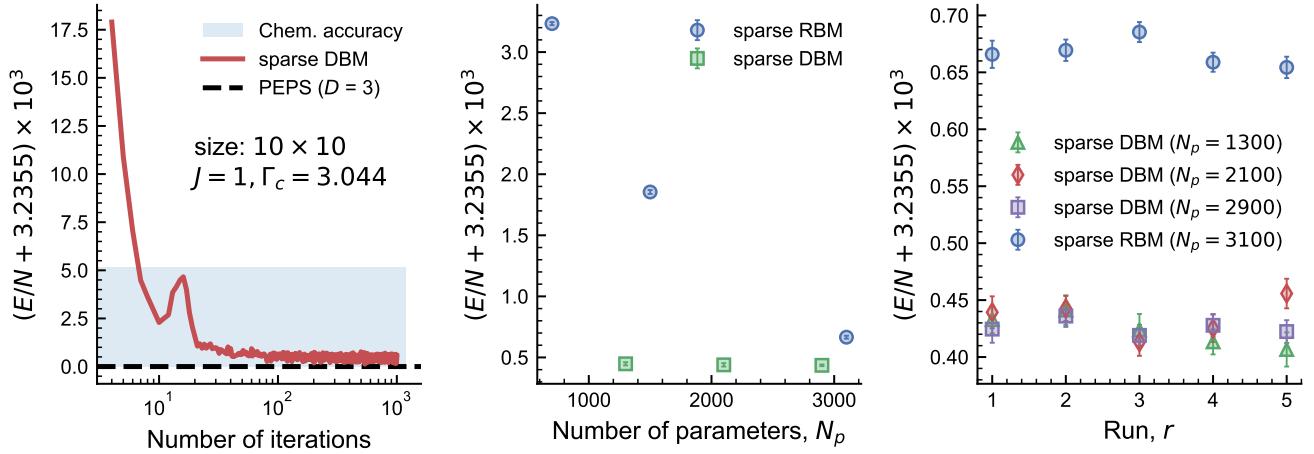
**Fig. 4**. **Deep Boltzmann machine results for the 2D transverse-field Ising model.** (a) Convergence of the sparse DBM using dual sampling on a $10 \times 10$ lattice ($J = 1$, $\Gamma_c = 3.044$). Supplementary Section E provides the hyperparameter choices. The variational energy converges toward a variational PEPS benchmark from Ref. [35]. A transient non-monotonic deviation appears early in training, consistent with the adaptive regularization dynamics of the SR solver during the high-learning-rate phase. (b) Final converged energy versus number of variational parameters $N_p$ for the sparse DBM (green squares) and sparse RBM (blue circles). Parameter count is varied by adjusting the connectivity radius $k$ (see Supplementary Section F). The most compact DBM ($N_p \approx 1300$) achieves lower energies than RBMs with more than twice as many parameters. (c) Final energy across 5 independent training runs with random initializations. The DBM consistently converges to lower energies than the RBM. Error bars in (b) and (c) represent standard error via blocking (50 bins); the y-axis shows residual energy scaled by $10^3$ relative to the offset $-3.2355$.
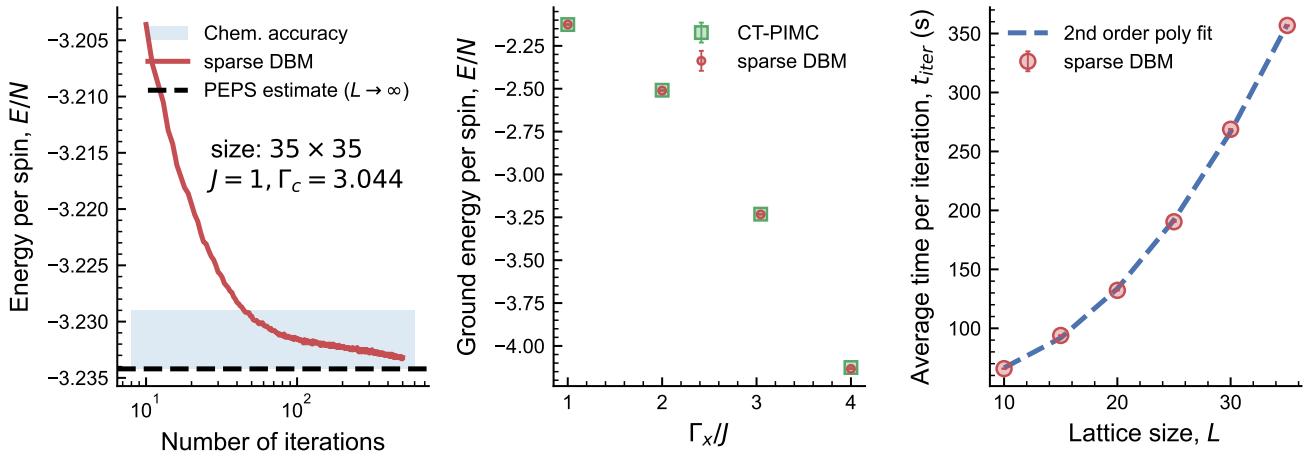


**Fig. 5**. **Scalable training of Deep Boltzmann Machines (DBMs) on GPU.** (a) Training convergence for a $35 \times 35$ lattice at the critical point $\Gamma_c/J \approx 3.044$. We train a sparse DBM with local connectivity radius $k = 2$ on both layers (13 connections per spin), resulting in $N_p = 35{,}525$ variational parameters. The variational energy (red line) converges to within chemical accuracy (blue shaded region) of a variational PEPS benchmark estimate at the thermodynamic limit ($L \to \infty$) from Ref. [35]. (b) Ground-state energy density across the transverse-field phase transition for a $35 \times 35$ lattice, comparing the sparse DBM (red circles) against CT-PIMC (green squares, same as in Fig. 2). The DBM captures the phase transition with high precision. (c) Empirical algorithmic scaling of the dual-sampling estimator on a single NVIDIA V100 GPU. The average wall-clock time per optimization iteration (red circles, averaged over 500 iterations; error bars indicate 95% confidence intervals) is shown as a function of the linear lattice dimension $L$. To strictly isolate algorithmic scaling from hardware parallelization effects, we use a single sequential MCMC chain rather than parallel independent chains. Under fixed sparsity and sampling budgets ($N_s$, $N_c$), the iteration time scales approximately as $t_{\text{iter}} \propto L^2$ (dashed quadratic fit), corresponding to linear dependence on the total number of spins $N = L^2$. This behavior is consistent with the simplified cost model for conventional processors in Supplementary Section G.

depth provides more efficient representations [51]. The improvement is modest here, likely because the TFIM has relatively simple entanglement structure; deeper gains may emerge for Hamiltonians with more complex correlations.

Finally, Fig. 4c shows the final energy across five independent training runs with random initializations. The sparse DBM consistently converges to lower energies, whereas the sparse RBM settles into higher-energy local minima. This reproducibility suggests that the additional layers help mediate interactions between distant lattice regions while preserving hardware-compatible sparsity. More generally, the dual-sampling framework is not restricted to RBM or layered DBM architectures, but applies to arbitrary Boltzmann-machine graphs, including networks

with unstructured, interlayer, and hierarchical connectivity. This generality enables systematic exploration of expressive graph topologies for neural quantum states beyond the constraints of analytically tractable models.

We now evaluate the *algorithmic* scaling of dual sampling under fixed sparsity and sampling budgets, independent of specialized sampling hardware. We use a conventional GPU to make sampling costs explicit and to measure how the dual-sampling estimator itself scales with system size on conventional processors.

This choice is conservative. Sparse, locally connected graphs are difficult to execute efficiently on GPUs due to irregular memory access and limited data reuse, unlike the dense linear algebra workloads for which GPUs are optimized [57–59]. Stable convergence and predictable scaling under these conditions therefore provide a stringent test of the dual-sampling estimator.

Sparsity is also a design choice that reduces parameter count, arithmetic operations, and memory traffic per iteration, and aligns naturally with the probabilistic hardware architecture [60, 61]. The results in Fig. 5 quantify this scaling and complement the FPGA implementation, which reduces the sampling latency for large systems.

To demonstrate the scalability of the dual sampling algorithm, we trained a sparse DBM on a $35 \times 35$ lattice using a single NVIDIA V100 GPU. As shown in Fig. 5a, the variational energy converges within chemical accuracy of the benchmark ground-state energy at the critical point ($\Gamma_c/J \approx 3.044$). This stable convergence reflects the variance reduction achieved by conditioning the inner-loop sampling on fixed visible configurations.

The accuracy of the trained DBM extends across the phase diagram. Fig. 5b compares the DBM ground-state energy density against CT-PIMC benchmarks, showing agreement throughout both the ferromagnetic and field-polarized phases.

Fig. 5c plots the average wall-clock time per iteration as a function of the linear lattice dimension $L$, revealing a quadratic dependence $t_{\mathrm{iter}} \propto L^2$ under fixed sparsity and sampling budgets. Since the total number of spins scales as $N = L^2$, this behavior is consistent with the simplified cost analysis presented in Supplementary Section G. Probabilistic hardware can remove the explicit system-size dependence from sampling under the on-chip mapping assumption. However, the overall training cost remains linear in $N$ (quadratic in $L$) due to stochastic reconfiguration updates performed on the host processor. Further reducing the iteration time would require accelerating these optimization routines as well, for example through dedicated hardware support, which we leave for future work.

## VII. Conclusion

We have demonstrated that probabilistic computers built from p-bits can extend neural quantum state simulations to system sizes beyond current software implementations. Mapping sparse Boltzmann machines onto a multi-FPGA cluster, we obtained ground-state energies within chemical

accuracy for the 2D transverse-field Ising model at criticality, reaching 6400 spins. We also introduced dual sampling, an algorithm that replaces intractable marginalization in deep Boltzmann machines with conditional sampling, enabling the training of deep architectures with improved parameter efficiency relative to shallow networks. These results suggest that spatial locality and asynchronous parallelism are the key computational resources for scalable variational Monte Carlo. The sparse connectivity that enables efficient hardware mapping also makes local energy evaluation embarrassingly parallel: each wavefunction ratio depends only on a bounded neighborhood, opening a path to computing $E_{\mathrm{loc}}$ entirely on-chip and eliminating the remaining data-transfer bottleneck. Similarly, the dual-sampling algorithm is naturally suited to hardware acceleration, since its inner loop consists entirely of conditional Gibbs sampling over sparse graphs.

As p-bit architectures mature from FPGA prototypes to dedicated CMOS circuits, these optimizations become increasingly attractive. Application-specific hardware could integrate sampling, local energy evaluation, and gradient accumulation on a single die, reducing both latency and energy consumption by orders of magnitude. Such a platform would make variational Monte Carlo practical for quantum systems far larger than those accessible today, positioning probabilistic computing as a scalable route to classically simulating quantum matter.

## Methods

### p-computing overview

p-computing is based on a network of p-bits $\{\sigma_i\}$, each fluctuating between two logical states ($\sigma_i \in \{-1, +1\}$). The interactions between p-bits govern their state updates according to the following equations [30]:

$$I_i = \sum_j W_{ij}\sigma_j + b_i \tag{8}$$

$$\sigma_i = \mathrm{sgn}\left(\tanh(\beta I_i) - r_{[-1,1]}\right) \tag{9}$$

where $W$, $b$, and $\beta$ represent the interconnection matrix, bias vector, and inverse temperature, respectively. $r_{[-1,1]}$ is a uniformly distributed random variable in the interval $[-1, 1]$. Together, equations (8) and (9) drive the network toward the Boltzmann distribution:

$$p(\{\sigma_i\}) = \frac{1}{Z}\exp\left(-\beta E(\{\sigma_i\})\right) \tag{10}$$

$$E(\{\sigma_i\}) = -\sum_{i<j} W_{ij}\sigma_i\sigma_j - \sum_i b_i\sigma_i \tag{11}$$

where $p(\{\sigma_i\})$ is the probability and $E(\{\sigma_i\})$ is the energy of the state $\{\sigma_i\}$, and $Z$ denotes the partition function.

## FPGA details

We mapped the physics-inspired, massively parallel *p*-computer architecture of Ref. [33] onto a Xilinx Alveo U250 accelerator card using graph coloring to maximize parallelism on the sparse instances. All arithmetic is fixed-point and uses s{6}{3} precision (1 sign, 6 integer, 3 fractional bits). Custom RTLs were developed to implement the algorithm based on the p-computing architecture and synthesized, placed and routed with Xilinx Vivado/Vitis tool chain. Weights and biases are converted to this fixed-point precision in MATLAB before being sent to the FPGA.

## DATA AVAILABILITY

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## CODE AVAILABILITY

The codes that support the findings of this study are available from the corresponding author upon reasonable request.

## Acknowledgments

## Author Contributions

SC, JHM and KYC conceived the study. SC and JP performed different parts of the simulations. NAA provided support for the multi-FPGA implementations. SC and JP wrote the initial draft of the manuscript with inputs from JHM and KYC. All authors contributed to improving the draft and participated in designing the experiments, analyzing the results, and editing the manuscript.

## Competing Interests

The authors declare no competing interests.

## References

[1] Matthias Troyer and Uwe-Jens Wiese. Computational Complexity and Fundamental Limitations to Fermionic Quantum Monte Carlo Simulations. *Physical review letters*, 94 (17):170201, 2005.

[2] Ulrich Schollwöck. The density-matrix renormalization group in the age of matrix product states. *Annals of Physics*, 326(1): 96–192, 2011.

[3] Román Orús. A practical introduction to tensor networks: Matrix product states and projected entangled pair states. *Annals of Physics*, 349:117–158, 2014.

[4] Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325):602–606, February 2017.

[5] Nicolas Le Roux and Yoshua Bengio. Representational power of restricted boltzmann machines and deep belief networks. *Neural Computation*, 20(6):1631–1649, 2008.

[6] Markus Schmitt, Marek M. Rams, Jacek Dziarmaga, Markus Heyl, and Wojciech H. Zurek. Quantum phase transition dynamics in the two-dimensional transverse-field ising model. *Science Advances*, 8(37):eabl6850, 2022.

[7] G. Fabiani, M. D. Bouman, and J. H. Mentink. Supermagnonic propagation in two-dimensional antiferromagnets. *Phys. Rev. Lett.*, 127:097202, Aug 2021.

[8] Yusuke Nomura and Masatoshi Imada. Dirac-type nodal spin liquid revealed by refined quantum many-body solver using neural-network wave function, correlation ratio, and level spectroscopy. *Phys. Rev. X*, 11:031034, Aug 2021.

[9] Kenny Choo, Giuseppe Carleo, Nicolas Regnault, and Titus Neupert. Symmetries and many-body excitations with neural-network quantum states. *Phys. Rev. Lett.*, 121:167204, Oct 2018.

[10] Mohamed Hibat-Allah, Martin Ganahl, Lauren E. Hayward, Roger G. Melko, and Juan Carrasquilla. Recurrent neural network wave functions. *Phys. Rev. Res.*, 2:023358, Jun 2020.

[11] Luciano Loris Viteritti, Riccardo Rende, and Federico Becca. Transformer variational wave functions for frustrated quantum spin systems. *Phys. Rev. Lett.*, 130:236401, Jun 2023.

[12] Yuan-Hang Zhang and Massimiliano Di Ventra. Transformer quantum state: A multipurpose model for quantum many-body problems. *Phys. Rev. B*, 107:075147, Feb 2023.

[13] Tianchen Zhao, Saibal De, Brian Chen, James Stokes, and Shravan Veerapaneni. Overcoming barriers to scalability in variational quantum monte carlo. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '21, New York, NY, USA, 2021. Association for Computing Machinery.

[14] Riccardo Rende, Luciano Loris Viteritti, Federico Becca, Antonello Scardicchio, Alessandro Laio, and Giuseppe Carleo. Foundation neural-networks quantum states as a unified ansatz

for multiple hamiltonians. *Nature Communications*, 16(1): 7213, Aug 2025.

[15] Penghua Chen, Bowen Yan, and Shawn X. Cui. Representing arbitrary ground states of the toric code by a restricted boltzmann machine. *Phys. Rev. B*, 111:045101, Jan 2025.

[16] Marc Machaczek, Lode Pollet, and Ke Liu. Neural quantum state study of fracton models. *SciPost Phys.*, 18:112, 2025.

[17] Giammarco Fabiani and Johan H. Mentink. Investigating ultrafast quantum magnetism with machine learning. *SciPost Phys.*, 7:004, 2019.

[18] Markus Schmitt and Markus Heyl. Quantum Many-Body Dynamics in Two Dimensions with Artificial Neural Networks. *Physical Review Letters*, 125(10):100503, September 2020. Publisher: American Physical Society.

[19] Filippo Vicentini, Damian Hofmann, Attila Szabó, Dian Wu, Christopher Roth, Clemens Giuliani, Gabriel Pescia, Jannes Nys, Vladimir Vargas-Calderón, Nikita Astrakhantsev, and Giuseppe Carleo. NetKet 3: Machine Learning Toolbox for Many-Body Quantum Systems. *SciPost Phys. Codebases*, page 7, 2022.

[20] Markus Schmitt and Moritz Reh. jVMC: Versatile and performant variational Monte Carlo leveraging automated differentiation and GPU acceleration. *SciPost Phys. Codebases*, page 2, 2022.

[21] Ao Chen and Markus Heyl. Empowering deep neural quantum states through efficient optimization. *Nature Physics*, 20(9): 1476–1481, Sep 2024. ISSN 1745-2481.

[22] M. Schuyler Moss, Roeland Wiersema, Mohamed Hibat-Allah, Juan Carrasquilla, and Roger G. Melko. Leveraging recurrence in neural network wavefunctions for large-scale simulations of heisenberg antiferromagnets on the triangular lattice. *Phys. Rev. B*, 112:134449, Oct 2025.

[23] M. Schuyler Moss, Roeland Wiersema, Mohamed Hibat-Allah, Juan Carrasquilla, and Roger G. Melko. Leveraging recurrence in neural network wavefunctions for large-scale simulations of heisenberg antiferromagnets on the square lattice. *Phys. Rev. B*, 112:134450, Oct 2025.

[24] Mingfan Li, Junshi Chen, Qian Xiao, Fei Wang, Qingcai Jiang, Xuncheng Zhao, Rongfen Lin, Hong An, Xiao Liang, and Lixin He. Bridging the gap between deep learning and frustrated quantum spin system for extreme-scale simulations on new generation of sunway supercomputer. *IEEE Transactions on Parallel and Distributed Systems*, 33(11):2846–2859, 2022.

[25] Xiao Liang, Mingfan Li, Qian Xiao, Junshi Chen, Chao Yang, Hong An, and Lixin He. Deep learning representations for quantum many-body systems on heterogeneous hardware. *Machine Learning: Science and Technology*, 4(1):015035, mar 2023.

[26] Hongtao Xu, Zibo Wu, Mingzhen Li, and Weile Jia. Large-scale neural network quantum states for ab initio quantum chemistry simulations on fugaku, 2025.

[27] Massimo Bortone, Yannic Rath, and George H. Booth. Impact of conditional modelling for a universal autoregressive quantum state. *Quantum*, 8:1245, Feb 2024.

[28] David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. A learning algorithm for boltzmann machines. *Cognitive Science*, 9(1):147–169, 1985.

[29] Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre G.R. Day, Clint Richardson, Charles K. Fisher, and David J. Schwab. A high-bias, low-variance introduction to machine learning for physicists. *Physics Reports*, 810:1–124, 2019. A high-bias, low-variance introduction to Machine Learning for physicists.

[30] Kerem Yunus Camsari, Rafatul Faria, Brian M Sutton, and Supriyo Datta. Stochastic p-bits for invertible logic. *Physical Review X*, 7(3):031014, 2017.

[31] Rutger J. L. F. Berns, Davi R. Rodrigues, Giovanni Finocchio, and Johan H. Mentink. Predicting sampling advantage of stochastic ising machines for quantum simulations, 2025.

[32] Pratik Brahma, Junghoon Han, Tamzid Razzaque, Saavan Patel, and Sayeef Salahuddin. Hardware acceleration of frustrated lattice systems using convolutional restricted boltzmann machine, 2025.

[33] Navid Anjum Aadit, Andrea Grimaldi, Mario Carpentieri, Luke Theogarajan, John M Martinis, Giovanni Finocchio, and Kerem Y Camsari. Massively parallel probabilistic computing with sparse Ising machines. *Nature Electronics*, 5(7):460–468, 2022.

[34] Shaila Niazi, Shuvro Chowdhury, Navid Anjum Aadit, Masoud Mohseni, Yao Qin, and Kerem Y. Camsari. Training deep Boltzmann networks with sparse Ising machines. *Nature Electronics*, pages 1–10, June 2024. Publisher: Nature Publishing Group.

[35] Gleb Fedorovich, Lukas Devos, Jutho Haegeman, Laurens Vanderstraeten, Frank Verstraete, and Atsushi Ueda. Finite-size scaling on the torus with periodic projected entangled-pair states. *Phys. Rev. B*, 111:165124, Apr 2025.

[36] Andrew D. King, Jack Raymond, Trevor Lanting, Sergei V. Isakov, Masoud Mohseni, Gabriel Poulin-Lamarre, Sara Ejtemaee, William Bernoudy, Isil Ozfidan, Anatoly Yu. Smirnov, Mauricio Reis, Fabio Altomare, Michael Babcock, Catia Baron, Andrew J. Berkley, Kelly Boothby, Paul I. Bunyk, Holly Christiani, Colin Enderud, Bram Evert, Richard Harris, Emile Hoskinson, Shuiyuan Huang, Kais Jooya, Ali Khodabandelou, Nicolas Ladizinsky, Ryan Li, P. Aaron Lott, Allison J. R. MacDonald, Danica Marsden, Gaelen Marsden, Teresa Medina, Reza Molavi, Richard Neufeld, Mana Norouzpour, Travis Oh, Igor Pavlov, Ilya Perminov, Thomas Prescott, Chris Rich, Yuki Sato, Benjamin Sheldan, George Sterling, Loren J. Swenson, Nicholas Tsai, Mark H. Volkmann, Jed D. Whittaker, Warren Wilkinson, Jason Yao, Hartmut Neven, Jeremy P. Hilton, Eric Ladizinsky, Mark W. Johnson, and Mohammad H. Amin. Scaling advantage over path-integral monte carlo in quantum simulation of geometrically frustrated magnets. *Nature Communications*, 12(1):1113, 02 2021.

[37] Subir Sachdev. *Quantum Phase Transitions*. Cambridge University Press, 2 edition, 2011.

[38] W. Marshall. Antiferromagnetism. *Proceedings of the Royal Society of London. A. Mathematical and Physical Sciences*, 232 (1188):48–68, 10 1955.

[39] Sergey Bravyi, David P. Divincenzo, Roberto Oliveira, and Barbara M. Terhal. The complexity of stoquastic local hamiltonian problems. *Quantum Info. Comput.*, 8(5):361–385, May 2008.

[40] Attila Szabó and Claudio Castelnovo. Neural network wave functions and the sign problem. *Phys. Rev. Res.*, 2:033075, Jul 2020.

[41] Giacomo Torlai, Guglielmo Mazzola, Juan Carrasquilla, Matthias Troyer, Roger Melko, and Giuseppe Carleo. Neural-network quantum state tomography. *Nature Physics*, 14(5): 447–450, May 2018.

[42] Shuvro Chowdhury, Navid Anjum Aadit, Andrea Grimaldi, Eleonora Raimondo, Atharva Raut, P. Aaron Lott, Johan H. Mentink, Marek M. Rams, Federico Ricci-Tersenghi, Massimo Chiappini, Luke S. Theogarajan, Tathagata Srimani, Giovanni Finocchio, Masoud Mohseni, and Kerem Y. Camsari. Pushing the boundary of quantum advantage in hard combinatorial optimization with probabilistic computers. *Nature Communications*, 16(1):9193, Oct 2025.

[43] Dong-Ling Deng, Xiaopeng Li, and S. Das Sarma. Machine learning topological states. *Phys. Rev. B*, 96:195145, Nov 2017.

[44] Yusuke Nomura, Andrew S. Darmawan, Youhei Yamaji, and Masatoshi Imada. Restricted boltzmann machine learning for solving strongly correlated quantum systems. *Phys. Rev. B*, 96: 205152, Nov 2017.

[45] David Blackman and Sebastiano Vigna. Scrambled linear pseudorandom number generators. *ACM Trans. Math. Softw.*, 47(4), September 2021.

[46] Shuvro Chowdhury, Kerem Y. Camsari, and Supriyo Datta. Accelerated quantum Monte Carlo with probabilistic computers. *Communications Physics*, 6(1):1–9, April 2023. Number: 1 Publisher: Nature Publishing Group.

[47] Sandro Sorella. Green function monte carlo with stochastic reconfiguration. *Phys. Rev. Lett.*, 80:4558–4561, May 1998.

[48] Eric Neuscamman, C. J. Umrigar, and Garnet Kin-Lic Chan. Optimizing large parameter sets in variational quantum monte carlo. *Phys. Rev. B*, 85:045103, Jan 2012.

[49] Srijan Nikhar, Sidharth Kannan, Navid Anjum Aadit, Shuvro Chowdhury, and Kerem Y Camsari. All-to-all reconfigurability with sparse and higher-order Ising machines. *Nature Communications*, 15(1):8977, 2024.

[50] Henk W. J. Blöte and Youjin Deng. Cluster monte carlo simulation of the transverse ising model. *Phys. Rev. E*, 66: 066110, Dec 2002.

[51] Xun Gao and Lu-Ming Duan. Efficient representation of quantum many-body states with deep neural networks. *Nature Communications*, 8(1):662, Sep 2017.

[52] Giuseppe Carleo, Yusuke Nomura, and Masatoshi Imada. Constructing exact representations of quantum many-body systems with deep neural networks. *Nature Communications*, 9 (1):5322, Dec 2018.

[53] Ruslan Salakhutdinov and Geoffrey E. Hinton. Deep boltzmann machines. In David A. Van Dyk and Max Welling, editors, *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics, AISTATS 2009, Clearwater Beach, Florida, USA, April 16-18, 2009*, volume 5 of *JMLR Proceedings*, pages 448–455. JMLR.org, 2009.

[54] Yusuke Nomura, Nobuyuki Yoshioka, and Franco Nori. Purifying Deep Boltzmann Machines for Thermal Quantum States. *Physical Review Letters*, 127(6):060601, August 2021. Publisher: American Physical Society.

[55] Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, August 2002.

[56] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, July 2006.

[57] Nathan Bell and Michael Garland. Implementing sparse matrix–vector multiplication on throughput-oriented processors. *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis (SC)*, 2009.

[58] Samuel Williams, Andrew Waterman, and David Patterson. Optimization of sparse matrix–vector multiplication on emerging multicore platforms. *Parallel Computing*, 35(3):178–194, 2009.

[59] Yangzihao Chen, Nadathur Satish, Sungpack Hong, Oluwasegun Oguntebi, and Kunle Olukotun. Gunrock: Gpu graph analytics. *ACM Transactions on Parallel Computing*, 4 (1):1–49, 2018.

[60] Mark Horowitz. 1.1 computing's energy problem (and what we can do about it). *IEEE International Solid-State Circuits Conference (ISSCC)*, pages 10–14, 2014.

[61] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019.

[62] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.

## Supplementary Information

# Probabilistic Computers for Neural Quantum States

Shuvro Chowdhury, Jasper Pieterse, Navid Anjum Aadit, Johan H. Mentink, and Kerem Y. Camsari

### A. Runtime Profiling and Bottleneck Analysis for 2D TFIM

To quantify the removal of the sampling bottleneck, we profiled the wall-clock execution time of the training loop across three hardware configurations for 2D TFIM with $L = 35$: a standard CPU implementation (11th Gen Intel Core i7-11700 @ 2.50 GHz, 96 GB RAM: also used in the hybrid setup), a MATLAB-based high-performance GPU implementation (NVIDIA GeForce RTX 4060 Ti, 16 GB VRAM), and our hybrid CPU-FPGA setup.

### B. Implementation Details of the Multi-FPGA System

We provide an extended view of the cluster architecture in Fig. S2, specifically detailing the mapping of partitioned subgraphs to hardware and the protocols used for inter-FPGA data exchange.

### C. Detailed Dual Sampling Algorithm

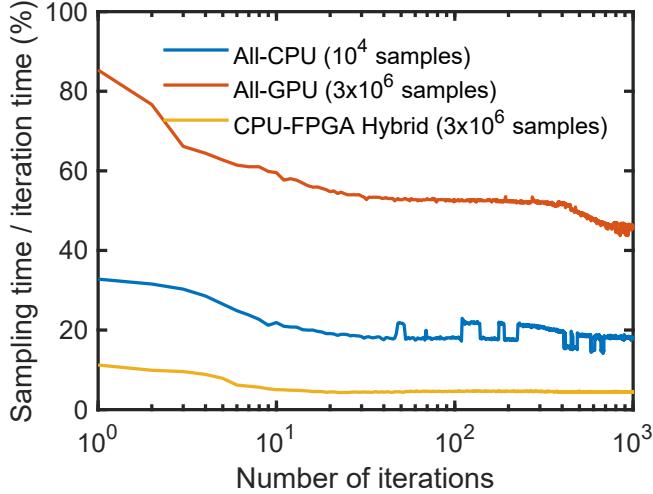Here, we provide a more elaborative version of the dual sampling algorithm.

**Fig. S1**. **Evolution of the computational bottleneck.** For 2D TFIM of size $L = 35$, the fraction of total iteration time spent on Monte Carlo sampling is plotted against training iterations for three architectures. The **All-GPU** setup (red, $3 \times 10^6$ samples) remains dominated by sampling costs ($> 50\%$), indicating a throughput bottleneck. The **All-CPU** baseline (blue) spends $\sim 20 - 30\%$ of its time on sampling despite processing $300\times$ fewer samples ($10^4$, for a manageable end-to-end training). In contrast, the **CPU-FPGA Hybrid** (yellow) performs the massive $3 \times 10^6$ sampling task in less than 5% of the total loop time, effectively rendering the sampling step instantaneous relative to the classical gradient accumulation.

### D.   Exactness of the dual-sampling estimator in the infinite-sample limit

Let us define a neural quantum state from a probabilistic DBM model $p_\theta(v)$ over spin configurations $v \in \{\pm 1\}^N$ by

$$\Psi_\theta(v) = \sqrt{p_\theta(v)}, \qquad \text{where} \quad p_\theta(v) = \sum_{h,d} p_\theta(v,h,d), \quad \text{and} \quad p_\theta(v,h,d) = \frac{1}{Z_\theta} e^{-E_\theta(v,h,d)} \tag{S.1}$$

so that $|\Psi_\theta(v)|^2 = p_\theta(v)$. Consider the transverse-field Ising Hamiltonian

$$H = -\sum_{\langle ij \rangle} J_{ij} \hat{\sigma}_i^z \hat{\sigma}_j^z - \Gamma_x \sum_{i=1}^N \hat{\sigma}_i^x$$

where the sum over $\langle ij \rangle$ includes only unique pairs connected by an edge in the underlying lattice. In the following, we will assume that

(1) in outer sampling, we can draw independent and identically distributed (i.i.d.) visible configurations from $p_\theta(v)$;

(2) in inner conditional sampling, for any fixed visible configuration $v$, we can sample the hidden/deep variables from their exact conditionals using the DBM and obtain an unbiased estimator of the ratio $p_\theta(v^{(i)})/p_\theta(v)$ for every single-spin flip $v^{(i)}$ (the configuration obtained from $v$ by flipping spin $i$);

For any state $\Psi_\theta$, the exact variational energy is defined as

$$E(\theta) = \frac{\langle \Psi_\theta | H | \Psi_\theta \rangle}{\langle \Psi_\theta | \Psi_\theta \rangle}$$

The usual variational Monte Carlo identity is

$$E(\theta) = \frac{\sum_v |\Psi_\theta(v)|^2 E_{\text{loc}}(v)}{\sum_v |\Psi_\theta(v)|^2} = \mathbb{E}_{v \sim |\Psi_\theta|^2} \left[ E_{\text{loc}}(v) \right], \tag{S.2}$$

where the *local energy* is

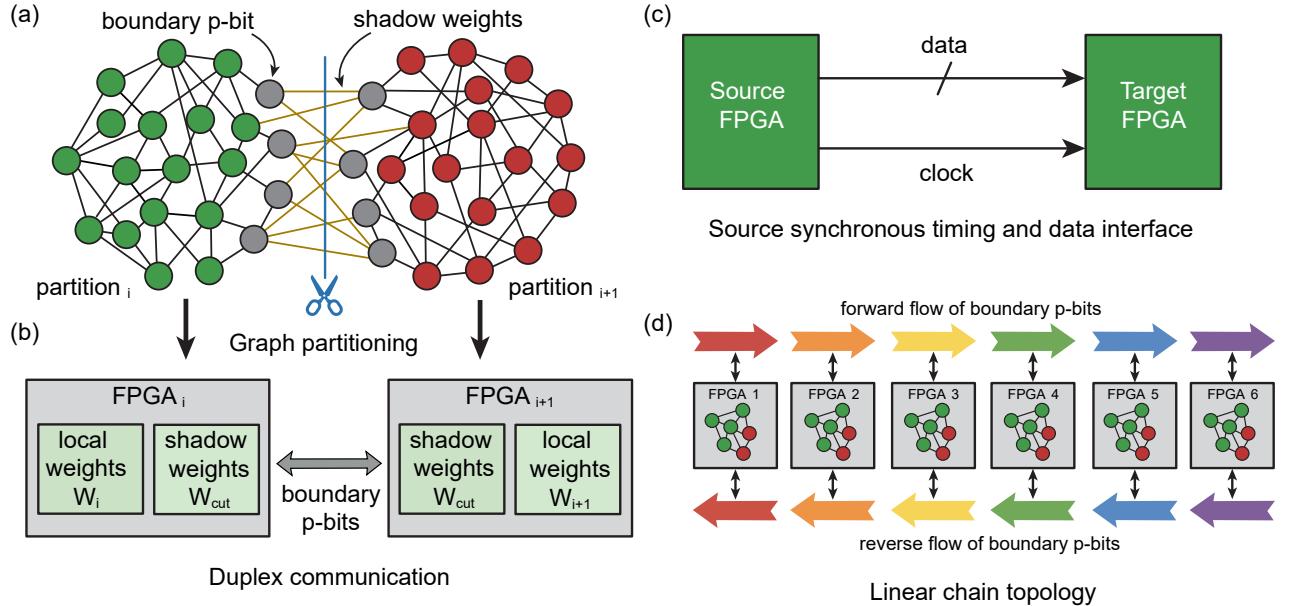$$E_{\text{loc}}(v) = \frac{(H\Psi_\theta)(v)}{\Psi_\theta(v)} \tag{S.3}$$

**Fig. S2**. **Multi-FPGA cluster architecture for large-scale NQS.** (a) A sparse Ising graph is partitioned into subgraphs using the METIS min-cut graph partitioner [62] and the subgraphs are mapped to neighboring FPGAs. (b) Each FPGA stores its local weights and duplicates cross-partition couplings as *shadow weights*. Only binary boundary p-bit states are exchanged over full-duplex FMC links (both directions) since the coupling matrix is symmetric. (c) Source-synchronous boundary interface: clock and data are forwarded together. Incoming boundary p-bits remain asynchronous to the destination FPGA clock domain and the forwarded clock is used only to align the received p-bit stream. (d) A 6-FPGA linear chain illustrating forward and reverse flow of boundary p-bit states over nearest-neighbor links.

For the TFIM Hamiltonian above, in the $\sigma^z$ basis, we have

$$E_{\text{loc}}(v) = -\sum_{\langle ij \rangle} J_{ij}\sigma_i^z\sigma_j^z - \Gamma_x \sum_{i=1}^{N} \frac{\Psi_\theta(v^{(i)})}{\Psi_\theta(v)} \tag{S.4}$$

Now, by construction of the NQS,

$$\Psi_\theta(v) = \sqrt{p_\theta(v)} \quad \Longrightarrow \quad \frac{\Psi_\theta(v^{(i)})}{\Psi_\theta(v)} = \sqrt{\frac{p_\theta(v^{(i)})}{p_\theta(v)}} = \sqrt{r_i(v)} \tag{S.5}$$

So if, for a fixed visible configuration $v$, we can obtain an unbiased estimate of $\frac{p_\theta(v^{(i)})}{p_\theta(v)}$, which is exactly what the inner (hidden/deep) sampling does in Algorithm S1, then applying the square root exactly recovers the amplitude ratio in Eq. (S.5). Let us denote the inner estimate (after averaging over $N_c$ conditional samples) by $\hat{r}_i(v)$. Under assumption (2) and letting $N_c \to \infty$, we can show that

$$\lim_{N_c \to \infty} \hat{r}_i(v) = \frac{p_\theta(v^{(i)})}{p_\theta(v)} \quad \Longrightarrow \quad \lim_{N_c \to \infty} \sqrt{\hat{r}_i(v)} = \sqrt{\frac{p_\theta(v^{(i)})}{p_\theta(v)}} = \frac{\Psi_\theta(v^{(i)})}{\Psi_\theta(v)}. \tag{S.6}$$

Plugging this into Eq. (S.4), we obtain an *inner-level* local-energy estimator

$$\lim_{N_c \to \infty} \hat{E}_{\text{loc}}(v) = -\sum_{\langle ij \rangle} J_{ij}\sigma_i^z\sigma_j^z - \Gamma_x \sum_{i=1}^{N} \left[ \lim_{N_C \to \infty} \sqrt{\hat{r}_i(v)} \right] = E_{\text{loc}}(v) \tag{S.7}$$

To see why Eq. (S.7) holds, consider a visible configuration $v$ and the configuration $v^{(i)}$ obtained by flipping spin $v_i$. For fixed hidden and deep spins $(h, d)$, the DBM defines a Boltzmann weight

$$p_\theta(v, h, d) = \frac{1}{Z_\theta} e^{-E_\theta(v,h,d)} \tag{S.8}$$

---

**Algorithm S1:** Algorithm for Machine Learning Quantum Hamiltonians using deep NQS

---

**Input:** Number of qubits ($N_Q$), parameters of the quantum Hamiltonian ($J, h, \Gamma$), number of samples per iterations ($N_s$), number of clamped samples ($N_c$), number of iterations ($N_{\text{iter}}$), learning rate schedule ($\eta$), optimization parameters

**Output:** Optimized parameters: biases ($b$), couplings ($W$), ground energy estimate ($E_G$)

1  Define DBM adjacency matrix, $A$.
2  Initialize biases: $b_i \leftarrow 0.01\,\mathcal{N}(0,1), \quad \forall i \in \{v, h, d\}$.
3  Initialize weights: $W_{ij} \leftarrow 0.01\,\mathcal{N}(0,1)$
4  Initialize p-bit states randomly from $\{-1, +1\}$.
5  **for** *each iteration $t = 1$ to $N_{\text{iter}}$*
6      **for** *each sample in $N_s$*
7        **for** *each spin in visible layer* **do in parallel**
8          Compute input $I_i = \sum_j W_{ij}\sigma_j + h_i$.
9          Update: $\sigma_i = \text{sgn}\left(\tanh(\beta I_i) - r_{[-1,1]}\right)$.
10       Clamp or fix the spins in the visible layer, $v$.
11       Compute diagonal contribution $H_{v,v}$.
12       $E_{\text{loc}}(v) \leftarrow H_{v,v}$.
13       $O_{v_i} \leftarrow O_{v_i} + 0.5\,v_i, \quad \forall i \in v$.
14       $O_{h_j}(v) \leftarrow 0, O_{d_k}(v) \leftarrow 0, O_{W_{ij}}(v) \leftarrow 0$.
15       **for** *each sample in $N_c$*
16         **for** *each spin in hidden layer* **do in parallel**
17           Compute input $I_i = \sum_j W_{ij}\sigma_j + h_i$.
18           Update: $\sigma_i = \text{sgn}\left(\tanh(\beta I_i) - r_{[-1,1]}\right)$.
19         $O_{h_j}(v) \leftarrow O_{h_j}(v) + 0.5\,h_j, \quad \forall j \in h$.
20         **for** *each spin in deep layer* **do in parallel**
21           Compute input $I_i = \sum_j W_{ij}\sigma_j + h_i$.
22           Update: $\sigma_i = \text{sgn}\left(\tanh(\beta I_i) - r_{[-1,1]}\right)$.
23         $O_{d_k}(v) \leftarrow O_{d_k}(v) + 0.5\,d_k, \quad \forall k \in d$.
24         Compute input $I_i, \quad \forall i \in v$.
25         $p_{\text{flip},i} \leftarrow p_{\text{flip},i} + \exp\left(-2I_i v_i\right), \quad \forall i \in v$.
26         $p_{\text{flip\_sq},i} \leftarrow p_{\text{flip\_sq},i} + \exp\left(-4I_i v_i\right), \quad \forall i \in v$.
27         $O_{W_{ij}}(v) \leftarrow O_{W_{ij}}(v) + 0.5\,\sigma_i\sigma_j, \ \forall i,j \in \{v,h,d\}$.
28       $O_{h_j}(v) \leftarrow O_{h_j}(v)/N_c, \quad \forall j \in h$.
29       $O_{d_k}(v) \leftarrow O_{d_k}(v)/N_c, \quad \forall j \in d$.
30       $O_{W_{ij}}(v) \leftarrow O_{W_{ij}}(v)/N_c, \ \forall i,j \in \{v,h,d\}$.
31       $p_{\text{flip},i} \leftarrow p_{\text{flip},i}/N_c, \quad \forall i \in v$.
32       $p_{\text{flip\_sq},i} \leftarrow p_{\text{flip\_sq},i}/N_c, \quad \forall i \in v$.
33       $\text{Var}_{\text{pop},i} \leftarrow (p_{\text{flip\_sq},i} - (p_{\text{flip},i})^2)/N_c \quad \forall i \in v$.
34       $\Delta_i \leftarrow \text{Var}_{\text{pop},i}/(8 p_{\text{flip},i}\sqrt{p_{\text{flip},i}}) \quad \forall i \in v$.
35       **for** *each spin in visible layer*
36         Get $v^{(i)}$ by flipping the spin.
37         $E_{\text{loc}}(v) \leftarrow E_{\text{loc}}(v) + H_{v,v^{(i)}}(\sqrt{p_{\text{flip},i}} + \Delta_i)$.
38     Use SR for gradients $\Delta b_i$ and $\Delta W_{ij}, \ \forall i,j \in \{v,h,d\}$.
39     Update biases: $b_i \leftarrow b_i - \eta(t)\Delta\,b_i, \ \forall i \in \{v,h,d\}$.
40     Update weights: $W_{ij} \leftarrow W_{ij} - \eta(t)\Delta W_{ij}, \ \forall i,j \in \{v,h,d\}$.
41  **return** Optimized $b$, $W$ and $E_G^{(N_{\text{iter}})}$.

---

and

$$\Delta\,E_i(v,h,d) = E_\theta(v^{(i)},h,d) - E_\theta(v,h,d) = 2I_i(v,h,d)v_i \tag{S.9}$$

so the ratio of joint probabilities is

$$\frac{p_\theta(v^{(i)},h,d)}{p_\theta(v,h,d)} = e^{-\Delta\,E_i(v,h,d)} = e^{-2I_i(v,h,d)v_i} \tag{S.10}$$

The marginal probability of a visible configuration is obtained by summing over hidden and deep variables:

$$p_\theta(v) = \sum_{h,d} p_\theta(v,h,d), \qquad p_\theta(v^{(i)}) = \sum_{h,d} p_\theta(v^{(i)}, h, d) \tag{S.11}$$

Using Eq. (S.10) in the second expression gives

$$p_\theta(v^{(i)}) = \sum_{h,d} p_\theta(v^{(i)}, h, d) = \sum_{h,d} p_\theta(v,h,d)\, e^{-\Delta\, E_i(v,h,d)} \tag{S.12}$$

Therefore,

$$\frac{p_\theta(v^{(i)})}{p_\theta(v)} = \frac{\sum_{h,d} p_\theta(v^{(i)}, h, d)}{\sum_{h,d} p_\theta(v,h,d)} = \frac{\sum_{h,d} p_\theta(v,h,d)\, e^{-\Delta\, E_i(v,h,d)}}{\sum_{h,d} p_\theta(v,h,d)} = \mathbb{E}_{(h,d)\sim p(h,d|v)}\left[e^{-\Delta\, E_i(v,h,d)}\right] = r_i(v) \tag{S.13}$$

which establishes Eq. (S.7). $r_i(v)$ is the conditional expectation, under the clamped distribution $p_\theta(h, d|v)$, of the joint Boltzmann ratio between $v^{(i)}$ and $v$. In our algorithm, this expectation is estimated by a Monte Carlo average over $N_c$ clamped samples,

$$\hat{r}_i(v) = \frac{1}{N_c} \sum_{k=1}^{N_c} e^{-\Delta\, E_i(v, h^{(k)}, d^{(k)})} \tag{S.14}$$

and by the law of large numbers $\hat{r}_i(v) \to r_i(v)$ as $N_c \to \infty$, so the estimator converges to the exact marginal ratio $\dfrac{p_\theta(v^{(i)})}{p_\theta(v)}$.

### E. Training Hyperparameters and Optimization Details

All Neural Quantum State (NQS) models presented in this work were trained using a custom MATLAB framework accelerated by CUDA-enabled GPUs. The optimization was performed using Stochastic Reconfiguration (SR) implemented via an iterative Conjugate Gradient (CG) solver to avoid explicitly constructing the dense Fisher Information Matrix (FIM).

#### 1. Hyperparameters

The specific hyperparameters used for the training of the Sparse Deep Boltzmann Machine (DBM) and Sparse RBM results shown in Figs. 4 and 5 are summarized in Table S1.

- **Dual Sampling Scheme:** For the Sparse DBM, the gradient estimation involves two distinct sampling populations. The *outer loop* uses $N_s = 10,000$ samples to estimate expectations over the visible layer $|\Psi(v)|^2$. The *inner loop* (conditional sampling) uses $N_c = 1,000$ clamped samples to estimate the gradients of the deep layers and the local energy ratios. For a fixed visible configuration, the same set of $N_c$ conditional samples is reused to estimate all single-spin-flip ratios, such that each ratio estimator is based on $N_c$ samples (not $N_c/N$).

- **Learning Rate Schedule:** For the algorithmic benchmarks in Fig. 4, we employed a time-dependent learning rate $\eta(t)$ following a *cosine decay schedule*, decreasing from $\eta_{\max} = 0.1$ to $\eta_{\min} = 10^{-5}$ over the course of training to ensure stable convergence. Conversely, for the hardware demonstrations in Fig. 2 and Fig. 3, a cosine decay schedule from $\eta_{\max} = 0.05$ to $\eta_{\min} = 0.01$ was employed.

- **Regularization:** To stabilize the inversion of the Fisher Information Matrix, we applied a diagonal shift regularization (Tikhonov regularization) $S \to S + \lambda I$. The shift parameter $\lambda$ was initialized at $\lambda_0 = 0.1$ and decayed adaptively during training to a minimum of $\lambda_{\min} = 10^{-4}$ using a decay factor $b_0 = 0.9$ per iteration ($\lambda_t = \max(\lambda_{\min}, \lambda_0 b_0^t)$).

The outer sampling budget $N_s$ and the inner conditional sampling budget $N_c$ control distinct sources of statistical error. The outer loop samples the full visible-spin distribution and determines the conditioning of the stochastic-reconfiguration update, while the inner loop estimates conditional expectations over auxiliary variables for a fixed visible configuration. Because the conditional distributions are lower-variance and strongly constrained, accurate estimation of wavefunction ratios can be achieved with substantially fewer conditional samples. In practice, we find that $N_c = 10^3$ is sufficient to stabilize training and achieve chemical accuracy for the systems studied here. We do not claim that this choice is asymptotically optimal. In general, the

hyperparameters listed in Table S1 denote fixed sampling budgets used throughout this work. We do not assume or claim that these quantities are independent of the system size. Rather, they are chosen empirically to ensure stable optimization and accurate energy estimation for the system sizes studied here. Determining how these sampling budgets must scale with system size to maintain a prescribed statistical accuracy would require a dedicated variance analysis and is beyond the scope of the present work.

## 2. Optimization Routine

The model parameters were updated using the Stochastic Reconfiguration (SR) method. Instead of inverting the curvature matrix $S$ directly (which scales as $\mathcal{O}(N_p^3)$), we solved the linear system $S \cdot \delta\theta = g$ using a Preconditioned Conjugate Gradient (PCG) solver.

- **Matrix-Free Implementation:** The solver utilizes implicit matrix-vector products to compute $S \cdot v$ without ever instantiating the full $S$ matrix in memory.

- **CG Tolerances:** The linear solver was configured with a relative tolerance of $10^{-4}$ and a maximum of 500 iterations per optimization step.

## 3. Final Evaluation

Following the training phase (1,000 iterations), the optimized model parameters were frozen. A final evaluation run was performed using a significantly larger sample size of $N_{\text{eval}} = 10^6$ to obtain the high-precision energy estimates and error bars reported in the main text figures.

TABLE S1. Summary of Training Hyperparameters

| Parameter | Symbol | Value | Description |
|---|---|---|---|
| Training Duration | $N_{\text{iter}}$ | 1000 | Total optimization steps |
| Visible Samples | $N_s$ | $10,000$ | Samples drawn from $|\Psi(v)|^2$ (Outer loop) |
| Clamped Samples | $N_c$ | $1,000$ | Conditional samples for Dual Sampling (Inner loop) |
| Initial Learning Rate | $\eta_{\max}$ | 0.1 | Maximum learning rate |
| Final Learning Rate | $\eta_{\min}$ | $10^{-5}$ | Minimum learning rate |
| SR Diagonal Shift | $\lambda_0$ | 0.1 | Initial regularization for Fisher Matrix |
| SR $\lambda$ decay factor | $b_0$ | 0.9 | Geometric decay rate for diagonal shift |
| CG Tolerance | $\epsilon_{\text{CG}}$ | $10^{-4}$ | Relative tolerance for linear solver |
| CG Max Iterations | $k_{\text{CG}}$ | 500 | Max iterations for linear solver |
| Evaluation Samples | $N_{\text{eval}}$ | $10^6$ | Samples used for final energy estimation |

## F. Parameter Counting and Sparse Connectivity Definitions

To quantify the representational efficiency of Deep Boltzmann Machines (DBMs) against Restricted Boltzmann Machines (RBMs), we constrain the network connectivity based on a geometric locality principle. Unlike standard "all-to-all" RBMs, which entail $\mathcal{O}(N^2)$ connections and are ill-suited for FPGA routing, we employ a spatially local connectivity mask.

### 1. Connectivity Metric

We define the connectivity between a neuron $i$ in layer $L$ at lattice coordinate $\mathbf{r}_i$ and a neuron $j$ in layer $L+1$ at $\mathbf{r}_j$ using the Euclidean distance on the periodic lattice (we assign to each hidden and deep layer a 2D geometry isomorphic to the visible lattice, such that every neuron has a defined spatial coordinate):

$$d(i,j) = \min_{\delta \in \mathbb{Z}^2} ||\mathbf{r}_i - \mathbf{r}_j + \mathbf{L} \cdot \delta||_2 \tag{S.15}$$

A synaptic connection $W_{ij}$ is non-zero if and only if $d(i,j) \leq k$, where $k$ is the tunable connectivity radius. We choose the Euclidean metric because it provides a simple and physically motivated notion of locality on the lattice, enabling sparse and distance-limited connectivity between neurons.

## 2. Parameter Enumeration

For the $10 \times 10$ lattice ($N = 100$ spins) used in Figure 4:

- **Sparse RBM:** Consists of one visible layer ($N_v = 100$) and one hidden layer ($N_h = 100$). The total parameters $N_p$ include $N_v + N_h$ biases and the number of active weights defined by $kN_v$.

- **Sparse DBM:** Consists of one visible layer ($N_v = 100$), one hidden layer ($N_h = 100$), and one deep layer ($N_d = 100$). The total parameters $N_p$ include $N_v + N_h + N_d$ biases and two sets of weights, constrained by connectivity radius $k_1$ (between visible-hidden) and $k_2$ (between hidden-deep).

Tables S2 and S3 details the exact parameter counts used for the sweep in Figure 4(b,c).

TABLE S2. Detailed parameter breakdown for Sparse RBM architectures evaluated in Figure 4 ($10 \times 10$ Lattice). The connectivity radius $k$ determines the number of active synapses per neuron.

| Radius ($k$) | Neighbors per Neuron | Weights | Total $N_p$ |
|---|---|---|---|
| 1 | 5 | 500 | 700 |
| 2 | 13 | 1300 | 1500 |
| 3 | 29 | 2900 | 3100 |

*Note:* "Neighbors per Neuron" includes the self-connection (same coordinate in adjacent layer) plus spatial neighbors. Total $N_p$ accounts for biases (200 for RBM).

TABLE S3. Detailed parameter breakdown for Sparse DBM architectures evaluated in Figure 4 ($10 \times 10$ Lattice). The connectivity radii ($k_1, k_2$) determine the number of active synapses per neuron.

| Radius ($k_1$) | Neighbors per Neuron | Radius ($k_2$) | Neighbors per Neuron | Weights | Total $N_p$ |
|---|---|---|---|---|---|
| 1 | 5 | 1 | 5 | 1000 | 1300 |
| 2 | 13 | 1 | 5 | 1800 | 2100 |
| 2 | 13 | 2 | 13 | 2600 | 2900 |

*Note:* "Neighbors per Neuron" includes the self-connection (same coordinate in adjacent layer) plus spatial neighbors. Total $N_p$ accounts for biases (300 for DBM).

## G. Simplified Computational Cost Model for Dual Sampling

We provide a simplified cost model for the dual-sampling algorithm used to estimate local energies and parameter derivatives for DBM neural quantum states in variational Monte Carlo. The goal of this section is to quantify how the wall-clock cost per iteration scales with model size at fixed sampling budgets and sparsity, and to demonstrate how probabilistic computing accelerates the sampling phase. We emphasize that the number of outer samples $N_s$ and inner conditional samples $N_c$ required to reach a prescribed statistical error can, in general, depend on the system size, Hamiltonian properties, sparsity and other parameters. Determining these dependencies requires a separate analysis and is outside the scope of this hardware-focused cost model.

In the following cost analysis, let $N_v$ denote the number of visible spins, $N_h$ the number of hidden units, and $N_d$ the number of deep units, with total number of stochastic units $N_u = N_v + N_h + N_d$. We assume a fixed sparse connectivity with bounded degree $k = O(1)$. Let $N_s$ be the number of outer Monte Carlo samples (visible configurations), and $N_c$ the number of conditional samples drawn in the inner loop for each visible configuration.

For each outer sample, one Gibbs sweep over the visible layer costs $O(N_v k)$. For each of $N_c$ conditional samples, a Gibbs sweep over auxiliary units costs $O((N_h + N_d)k)$ and evaluating all single-spin-flip ratios reusing the same conditional sample costs $O(N_v k)$ (exploiting the locality of the energy function). Thus, the sampling work per outer sample is

$$O(N_v k + N_c [(N_h + N_d)k + N_v k]), \tag{S.16}$$

and the total sampling work per optimization step is

$$O(N_s N_c N_u k). \tag{S.17}$$

For bounded degree $k = O(1)$, the total work scales as the product of the sampling budget and the system size, $O(N_s N_c N_u)$. Crucially, reusing each conditional sample to evaluate all single-spin-flip ratios avoids the quadratic $O(N_s N_c N_v N_u)$ work (or $O(N_s N_c N_u^2)$ when $N_v = O(N_u)$) that would arise if conditional samples were generated independently for each of the $N_v$ ratios.

On probabilistic hardware (p-computers) implemented on FPGAs, stochastic units can be updated in parallel subject to a fixed update schedule (e.g., graph coloring). Let $C$ denote the number of update phases, which depends on the local connectivity pattern but is independent of $N_u$ for bounded degree graphs. Assuming the full sparse Boltzmann network fits within hardware resources, the wall-clock time for one Gibbs sweep is $O(C)$. Therefore, the sampling *latency* per optimization step scales as

$$O(N_s N_c C). \tag{S.18}$$

Comparing Eq. (S.18) to Eq. (S.17), probabilistic hardware removes the explicit $N_u$ factor associated with sequential sampling updates on conventional processors. This is the key distinction between conventional processors and probabilistic hardware: for sparse models with bounded degree and fixed sampling budgets $(N_s, N_c)$, the sampling *latency* on probabilistic hardware scales as $O(N_s N_c C)$, where $C$ depends only on local connectivity and is independent of $N_u$ under the on-chip mapping assumption. Hardware resources (number of p-bits and synapses) scale approximately linearly with $N_u$, but the sampling rate remains constant with system size.

Finally, parameter updates are computed using stochastic reconfiguration (SR) with a matrix-free conjugate gradient (CG) solver. Let $N_p$ be the number of variational parameters and $K_{\mathrm{CG}}$ the number of CG iterations, which may in general vary with $N_u$. Each CG iteration requires sample-averaged matrix-vector products with cost $O(N_s N_p)$, giving

$$O(K_{\mathrm{CG}} N_s N_p) \tag{S.19}$$

per optimization step. For bounded-degree sparse graphs, $N_p$ scales linearly with system size, i.e., $N_p = O(N_u)$. If translational symmetry is enforced via parameter tying, the effective $N_p$ can be reduced further, but we treat the untied case here.