# Reinforcement Learning for Many-Body Ground-State Preparation Inspired by Counterdiabatic Driving

Jiahao Yao, 1,\* Lin Lin, 1,2,3 and Marin Bukov<br/>4,5,†

<sup>1</sup>Department of Mathematics, University of California, Berkeley, California 94720, USA
<sup>2</sup>Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA
<sup>3</sup>Challenge Institute for Quantum Computation, University of California, Berkeley, California 94720, USA
<sup>4</sup>Department of Physics, University of California, Berkeley, California 94720, USA
<sup>5</sup>Department of Physics, St. Kliment Ohridski University of Sofia,

5 James Bourchier Boulevard, 1164 Sofia, Bulgaria
(Dated: October 5, 2021)

The quantum alternating operator ansatz (QAOA) is a prominent example of variational quantum algorithms. We propose a generalized QAOA called CD-QAOA, which is inspired by the counterdiabatic (CD) driving procedure, designed for quantum many-body systems and optimized using a reinforcement learning (RL) approach. The resulting hybrid control algorithm proves versatile in preparing the ground state of quantum-chaotic many-body spin chains by minimizing the energy. We show that using terms occurring in the adiabatic gauge potential as generators of additional control unitaries, it is possible to achieve fast high-fidelity many-body control away from the adiabatic regime. While each unitary retains the conventional QAOA-intrinsic continuous control degree of freedom such as the time duration, we consider the order of the multiple available unitaries appearing in the control sequence as an additional discrete optimization problem. Endowing the policy gradient algorithm with an autoregressive deep learning architecture to capture causality, we train the RL agent to construct optimal sequences of unitaries. The algorithm has no access to the quantum state, and we find that the protocol learned on small systems may generalize to larger systems. By scanning a range of protocol durations, we present numerical evidence for a finite quantum speed limit in the nonintegrable mixed-field spin-1/2 Ising and Lipkin-Meshkov-Glick models, and for the suitability to prepare ground states of the spin-1 Heisenberg chain in the long-range and topologically ordered parameter regimes. This work paves the way to incorporate recent success from deep learning for the purpose of quantum many-body control.

### I. INTRODUCTION

The ability to prepare a quantum many-body system in its ground state is an important milestone in the quest for understanding and identifying novel collective quantum phenomena. The degree to which ground states can be confidently prepared in present-day quantum simulators, delineates the limits of our capabilities to investigate the properties of new materials or molecules, and to propose innovative technological applications based on quantum effects, such as high-temperature superconductors and superfluids, magnetic field sensors, topological quantum computers, or synthetic molecules.

Quantum simulators, such as ultracold and Rydberg atoms [1, 2], trapped ions [3–6], nitrogen vacancy centers [7–9], and superconducting qubits [6, 10], all require the development of state preparation schemes via realtime dynamical processes. Despite their high level of controllability, finding short protocols to prepare strongly-correlated ground states under platform-specific constraints, is a challenging problem in AMO-based quantum simulation platforms, due to the exponentially large Hilbert space dimensions of quantum many-body systems. On this background, speed-efficient protocols

also become progressively more important for near-term quantum computing devices [11], where simulation errors grow with the protocol duration due to imperfections in the implementation of the basic gate operations.

Developing versatile methods for ground state preparation will enable quantum simulators to investigate hitherto unexplored quantum phases of matter, and determine the behavior of order parameters, correlation lengths and critical exponents. Theoretically, although an exact mathematical expression for the ground state might be known in some models, it remains still largely unclear how to prepare it in a unitary dynamical process. In generic models, the lack of closed-form analytical solutions motivates the use of numerical algorithms. Prominent examples for quantum state preparation include established quantum control algorithms, such as GRAPE [12] and CRAB [13], and variational quantum eigensolvers (VQE) [14], such as the quantum approximate optimization ansatz (QAOA) [15].

In this study, we present a novel hybrid reinforcement learning (RL)/optimal control algorithm based on an autoregressive deep learning architecture. We improve the current state-of-the-art for digital quantum control techniques by enhancing the capabilities to find optimal protocols that prepare the ground state of quantum many-body systems. The emerging versatile algorithm combines discrete and continuous control parameters to achieve maximum flexibility in its applicability to a num-

<sup>\*</sup> jiahaoyao@berkeley.edu

<sup>†</sup> mgbukov@phys.uni-sofia.bg

ber of different models.

To cope with the complexity of preparing ordered states in quantum many-body systems, we introduce a novel ansatz inspired by variational gauge potentials and counter-diabatic (CD) driving [16–19]. This allows us to excite the system away from equilibrium in a controllable manner to find short high-fidelity protocols away from the adiabatic regime. We demonstrate that combining features of CD driving with the digital simulation character of conventional QAOA yields superior performance over a wide range of protocol durations and physical models. Compared to the standard counter-diabatic driving algorithms, CD-QAOA represents a more flexible ansatz which allows us to take into account (i) experimental constraints, such as drift terms that cannot be switched off, and (ii) control degrees of freedom not present in CD driving; (iii) CD-QAOA is not tied to a drive protocol which obeys specific boundary conditions (such as vanishing protocol speed). Unlike continuous CD driving, CD-QAOA offers a simple and easy-to-apply variational ansatz without reference to the exact ground state of the system, paving the way for versatile digital quantum control.

In particular, our RL agent constructs unitary protocols that transfer the population into the ground state of three nonintegrable spin models (spin-1/2 and spin-1 mixed-field Ising chains, and the anisotropic spin-1 Heisenberg chain) which feature long-range and topological order, and the integrable Lipkin-Meshkov-Glick (LMG) model which allows us to present simulations for a large number of particles. We show numerical evidence for the existence of a finite quantum speed limit in the nonintegrable mixed-field spin-1/2 Ising model: an almost perfect system-size scaling indicates that this behavior persists in the thermodynamic limit. Our RL agent has no access to unmeasurable quantum states which grow exponentially with the number of degrees of freedom in the system: this allows the protocols we find to generalize across a number of system sizes [for the spin-1/2 mixed-field Ising model, opening up the door to apply ideas of transfer learning to quantum manybody control. Finally, we demonstrate that the CD-QAOA ansatz has direct practical implications in digital quantum control: it leads to much shorter circuit depths while simultaneously improves the fidelity of the prepared state, which can be utilized to reduce detrimental errors in modern quantum computers.

# II. GENERALIZED CONTINUOUS-DISCRETE QUANTUM APPROXIMATE OPTIMIZATION ANSATZ

To prepare many-body quantum states, we seek a unitary process U which brings the system from a given initial state  $|\psi_i\rangle$  to the ground state  $|\psi_{\rm GS}\rangle$  of the Hamiltonian H (which we call the target state  $|\psi_*\rangle$ ). Typically, Hamiltonians can be decomposed as a sum of two

non-commuting parts  $H = H_1 + H_2$ , e.g. the kinetic and interaction energy. We want to construct

$$U(\{\alpha_j\}_{j=1}^q, \tau) = \prod_{j=1}^q U(\alpha_j, \tau_j)$$
 (1)

from a sequence  $\tau$  of q consecutive unitaries (or their generators)  $\tau_j$  chosen from a set  $\mathcal{A}$ , with  $\tau_j \neq \tau_{j+1}$ . Each  $U(\alpha_j, \tau_j)$  is parametrized by a continuous degree of freedom  $\alpha_j$  (e.g. time or rotation angle), i.e.  $U(\alpha_j, \tau_j) = \exp(-i\alpha_j\tau_j)$ . We formulate state preparation as an optimization problem which consists of determining (i) the sequence  $\tau$ , and (ii) the values of the variational parameters  $\alpha_j$ , such that  $U|\psi_i\rangle \approx |\psi_{\rm GS}\rangle$ .

Our goal is to prepare the ground state of a Hamiltonian H, without having access to the ground state itself. Therefore, we use energy as a cost function

$$E(\{\alpha_j\}_{j=1}^q, \tau) = \langle \psi_i | U^{\dagger}(\{\alpha_j\}_{j=1}^q, \tau) H U(\{\alpha_j\}_{j=1}^q, \tau) | \psi_i \rangle,$$
(2)

or energy-density E/N which has a well-behaved limit when increasing the number of particles N [20]. We denote the ground state energy by  $E_{\rm GS} = \langle \psi_{\rm GS} | H | \psi_{\rm GS} \rangle$ .

Note that conventional QAOA is recovered as a special case where one only considers two unitaries  $U_i =$  $U(\alpha_i, H_i) = \exp(-i\alpha_i H_i), j = 1, 2, \text{ and } \tau \text{ is one of the two}$ alternating sequences. Whenever nested commutators of  $H_i$  span the entire Lie algebra which generates transport on the complex projective space associated with the Hilbert space  $\mathcal{H}$  of the system, applying QAOA is already enough to prepare any state, provided that the underlying circuit depth q is sufficiently large, and the optimal  $\alpha_i$  can be found [21]. While true in theory, this is often impractical, since (i) it requires access to in principle unbounded durations, (ii) it increases the number of optimization parameters  $\alpha_i$ , and – with it – the probability to get stuck in a local minimum of the control landscape, and (iii) the condition that nested commutators of  $H_i$ span the entire Lie algebra is generally not satisfied for the  $H_i$ 's of interest in quantum many-body physics due to, e.g., symmetry constraints.

The generalized QAOA ansatz [Eq. (1)] allows us to utilize a larger set of unitaries  $\mathcal{A}$  to construct the optimal sequence and to reduce the circuit depth q. Inspired by counter-diabatic (CD) driving, we find that a particularly suitable choice in the context of quantum manybody state manipulation, is given by the operators in the adiabatic gauge potential series [Sec. III]. Therefore, we call the resulting algorithm CD-QAOA. A different ansatz using more than two unitaries was considered in Ref. [22].

Compared to conventional QAOA, CD-QAOA introduces a discrete high-level optimization to find the optimal protocol sequence  $\tau$ . The combined optimization landscape can be particularly difficult to navigate, due to the existence of so-called barren plateaus where exponentially many directions have vanishing gradients [23–26].

Additionally, the total number of all allowed protocol sequences,  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1}$  [27], scales exponentially with the number of unitaries q, and presents a challenging discrete combinatorial optimization problem per se; indeed, state preparation, formulated as optimization, can feature a glassy landscape [28, 29] [App. D]. However, overcoming these potential difficulties is associated with a potential gain: CD-QAOA allows retaining the flexibility offered by continuous optimization, while increasing the number of independent discrete control degrees of freedom to  $|\mathcal{A}|$ ; this enables us to reach larger parts of the Hilbert space in shorter durations, and with a smaller circuit depth, as compared to conventional QAOA.

Thus, we formulate ground state preparation as a two-level optimization scheme [30]. (1) Low-level optimization: given a fixed sequence  $\tau$ , we find the optimal values of  $\alpha_i$  using a continuous optimization solver, e.g. SLSQP [31] [App. B]. To cope with the associated rugged optimization landscape [App. D], we run multiple realizations of random initial conditions and post-select the values which yield minimum energy. This continuous optimization problem is also present in conventional QAOA. (2) High-level optimization: in addition to the low-level optimization, we also perform a discrete optimization for the sequence  $\tau$  itself, to determine the optimal order in which unitaries from the set A should occur. To tackle this combinatorial problem, we formulate the high-level optimization as a reinforcement learning (RL) problem. We learn the optimal protocol using Proximal Policy Optimization, a variant of policy gradient. The policy is parameterized by a deep autoregressive network, which allows choosing the control unitaries  $U(\alpha_i, \tau_i)$  sequentially. In practice, we sample a batch of sequences from the policy, evaluate the energy of each sequence in the low-level optimization, and apply policy gradient to update the parameters of the policy. This two-level optimization procedure is repeated in a number of training episodes until convergence [App. A].

# III. VARIATIONAL STATE PREPARATION INSPIRED BY COUNTER-DIABATIC DRIVING

A natural question arises as to how to choose the set  $\mathcal{A}$  of unitaries for the generalized discrete-continuous QAOA ansatz. One possibility is to consider a set of universal elementary quantum gates, e.g., in the context of a quantum computer [32, 33], and in this case  $\alpha_j$  are angles of rotation. We leave this exciting possibility for a future study, and focus here on many-body ground state preparation instead.

The complexity of many-body systems motivates the use of a physics-informed approach to defining the control unitaries in  $\mathcal{A}$ . Suppose we initialize the system in the ground state of the parent Hamiltonian  $H(\lambda=0)$ ; we target the ground state of  $H(\lambda=1)$ , seeking the functional form of a time-dependent protocol  $\lambda(t)$ . If the instantaneous ground state of  $H(\lambda)$  remains gapped during the

evolution, the adiabatic theorem guarantees the existence of a solution  $\lambda(t)$ ,  $t \in [0,T]$ , provided T is large compared to the smallest inverse gap along the adiabatic trajectory. However, when the gap is known to close (e.g. across a phase transition), or when the state population transfer has to be done fast, adiabatic state preparation fails.

Compared to the adiabatic paradigm, gauge potentials provide additional control directions in Hilbert space which enable paths that non-adiabatically lead to the target state in a short time. In many-body systems, it is not known in general how to determine the exact gauge potential required for CD driving. However, it is possible to define variational approximations [34, 35] using an operator-valued series expansion [App. E] similar to a Schrieffer-Wolff transformation [36], or Shortcuts to Adiabaticity methods [33, 37]. Nonetheless, recent numerical simulations suggest that the exact gauge potential in generic many-body systems is a non-local operator [34, 38] which renders the series expansion asymptotic.

For these reasons, here we consider the constituent terms to every order of the variational gauge potential series,  $H_j$ , independently, and use them to generate the set of unitaries  $\mathcal{A} = \{\mathrm{e}^{-i\alpha_j H_j}\}$  for CD-QAOA [39]. We emphasize that our CD-QAOA ansatz is not designed to approximate the gauge potential itself, as opposed to Ref. [40], yet it yields similar benefits w.r.t. preparing the target state. In Sec. V we compare directly our approach with the variational gauge potential ansatz from Ref. [34].

Since CD-QAOA is a generalization of QAOA aimed to be useful in practice, we need to ensure the accessibility of the control terms  $H_j$ . Because they appear in the first few orders of the gauge potential series,  $H_j$  are (sums of) local many-body operators [cf. App. E]. Thus, in principle, there is no physical obstruction to emulate them in the lab, although this depends on the details of the experimental platform (especially for the interaction terms). Additionally, in the context of many-body systems where energy is extensive, in order to guarantee that we do not tap into a source of infinite energy, we constrain the norm of the generators  $\alpha_j H_j$ : we view  $\alpha_j \geq 0$  as time durations, and fix  $\sum_{j=1}^q \alpha_j = T$ , with T the total protocol duration. This keeps  $\alpha_j$  on the same order of magnitude as the coupling constants in the parent Hamiltonian whose ground state we want to prepare.

# IV. MANY-BODY GROUND STATE PREPARATION

We consider four non-integrable many-body systems of increasing complexity: the spin-1/2 and spin-1 mixed-field Ising models, the spin-1 Heisenberg model, and the integrable Lipkin-Meshkov-Glick (LMG) model where a large number of degrees of freedom is accessible in a classical simulation. The goal of the RL agent is to prepare their ordered ground states, starting from a product state. To generate training data, we compute numerically

short-hand notation	spin operator $H_j$
X	$\sum_i S_i^x$
Z	$\sum_i S_i^z$
Z Z	$\sum_{z} S_{i}^{z} S_{i+1}^{z}$
Z Z+Z	$\sum_{i} JS_{i}^{z}S_{i+1}^{z} + h_{z}S_{i}^{z}$
Y	$\sum_{i} S_{i}^{y}$
XY	$\sum_{i} S_{i}^{x} S_{i}^{y} + S_{i}^{y} S_{i}^{x}$
YZ	$\sum_{i} S_i^y S_i^z + S_i^z S_i^y$
X Y	$\sum_{i} S_{i}^{x} S_{i+1}^{y} + S_{i}^{y} S_{i+1}^{x} $ $\sum_{i} S_{i}^{y} S_{i+1}^{z} + S_{i}^{z} S_{i+1}^{y}$
$Y Z \ X Y\!-\!XY$	$ \sum_{i} S_{i} S_{i+1} + S_{i} S_{i+1}  \sum_{i} [S_{i+1}^{x} - aS_{i}^{x}] S_{i}^{y} + [S_{i+1}^{y} - aS_{i}^{y}] S_{i}^{x} $
Y Z-YZ	$\sum_{i} [S_{i+1}^z - bS_i^z] S_i^y + [S_{i+1}^y - bS_i^y] S_i^z$
$\hat{XY}$	$\frac{1}{N} \sum_{i,j} S_i^x S_j^y + S_i^y S_j^x$
$\hat{ZY}$	$\frac{\frac{1}{N}\sum_{i,j}\left(S_i^z + \frac{I}{2}\right)S_j^y + S_i^y\left(S_j^z + \frac{I}{2}\right)}{\frac{1}{N}\sum_{i,j}\left(S_i^z + \frac{I}{2}\right)S_j^y + S_i^y\left(S_j^z + \frac{I}{2}\right)}$

TABLE I. Short-hand notation for the generators  $H_j$  used to construct the set of unitaries  $\mathcal{A} = \{e^{-i\alpha_j H_j}\}_{j=1}^{|\mathcal{A}|}$  in CD-QAOA. The | indicates operators acting on neighboring sites. Terms from the variational gauge potential series are shown in the lower group [cf. App. E for the derivation].

the exact time evolution of the system. We apply CD-QAOA using a set of unitaries built from the terms in the series expansion for the variational gauge potential. To determine the allowed terms in the gauge potential series, cf. Table I (lower group), we consider the minimal set of symmetries shared by the Hamiltonian and the initial and target states [App. E].

#### A. Mixed-Field Spin-1/2 Ising Chain

Consider first the antiferromagnetic mixed-field spin- 1/2 Ising chain of N lattice sites

$$H = H_1 + H_2,$$

$$H_1 = \sum_{j=1}^{N} J S_{j+1}^z S_j^z + h_z S_j^z, \quad H_2 = \sum_{j=1}^{N} h_x S_j^x,$$
(3)

where  $[S_i^{\alpha}, S_j^{\beta}] = \delta_{ij} \varepsilon^{\alpha\beta\gamma} S_j^{\gamma}$  are the spin-1/2 operators. We use periodic boundary conditions and work in the zero momentum sector of positive parity. In the following, J=1 sets the energy unit, and  $h_z/J=0.809$  and  $h_x/J=0.9045$ . We initialize the system in the z-polarized product state  $|\psi_i\rangle=|\uparrow\cdots\uparrow\rangle$ , and we want to prepare the ground state of H in a short time T, i.e., away from the adiabatic regime. We verified that similar results can be obtained starting from  $|\downarrow\cdots\downarrow\rangle$ .

To acquire an intuitive understanding of the advantages brought by the gauge potential ansatz, consider first the non-interacting system at J=0, for which the control problem reduces to a single spin. Both the initial and target states lie in the xz-plane of the Bloch sphere, and hence the shortest unit-fidelity protocol generates a rotation about the y-axis. In conventional QAOA, one would construct a y-rotation out of the X and Z terms

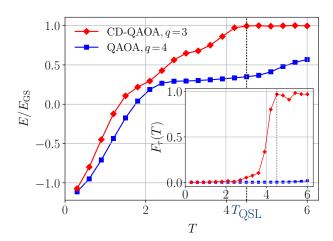


FIG. 1. Spin-1/2 Ising model: energy minimization and the corresponding many-body fidelity [inset] against protocol duration T obtained using conventional QAOA (blue squares) and CD-QAOA (red diamonds) with circuit depths p=q/2=2 and q=3, respectively. The dotted vertical line marks the quantum speed limit  $T_{\rm QSL}$ . CD-QAOA outperforms conventional QAOA. The initial and target states are  $|\psi_i\rangle=|\uparrow\cdots\uparrow\rangle$  and  $|\psi_*\rangle=|\psi_{\rm GS}(H)\rangle$  for  $h_z/J=0.809$  and  $h_x/J=0.9045$ . The alternating unitaries for conventional QAOA are generated by  $\mathcal{A}_{\rm QAOA}=\{Z|Z+Z,X\}$ [cf. Eq. (3)]; for CD-QAOA, we extend this set using adiabatic gauge potential terms to  $\mathcal{A}_{\rm CD-QAOA}=\{Z|Z+Z,X;Y,X|Y,Y|Z\}$ . The cardinality of the CD-QAOA sequence space is  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1}=80$ . The number of spins is N=18 with a Hilbert space size of  $\dim(\mathcal{H})=7685$ .

[cf. Table I] present in the Hamiltonian. For a single spin, this construction is always possible due to the Euler angle representation of SU(2), but for the interacting spin chain this is no longer the case. The role of the gauge potential Y is to 'unlock' precisely this geodesic in parameter space, and make it accessible as a dynamical process. This allows preparing the target state faster, compared to the original X, Z control setup. In the language of variational optimization, an accessible Y term includes the shortest-distance protocol into the variational manifold, and the RL agent easily finds the exact solution [App. F 1].

For the interacting system, J>0, applying conventional QAOA using the two gates  $U_j=\mathrm{e}^{-i\alpha_jH_j}$  with  $H_1=Z|Z+Z$  and  $H_2=X$  is straightforward, but it does not yield a high-fidelity protocol [Fig. 1 (blue squares)]. It was recently reported that much better energies can be obtained, using a three-step QAOA which consists of the three terms in the Hamiltonian (3), Z|Z, X, and Z, applied in a fixed order [41]; invoking again an Euler angle argument provides an explanation: the X and Z terms effectively generate the Y gauge potential term.

In stark contrast to conventional QAOA, adding just the zero-order term  $H_3 = Y$  from the gauge potential series [App. E3], we find that CD-QAOA already gives

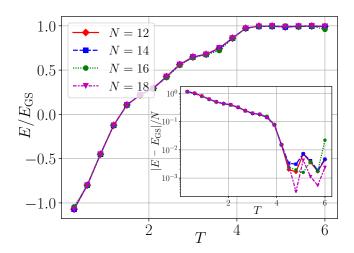


FIG. 2. Spin-1/2 Ising model: energy minimization and the corresponding mean absolute error [inset, log scale] against protocol duration T for different system sizes using CD-QAOA with circuit depths q=3. system-size scaling of the variational energy density suggests the results hold for larger systems. For the number of spins of N=12,14,16,18, the Hilbert space sizes are  $\dim(\mathcal{H})=224,687,2250,7685$  respectively. The model parameters are the same as in Fig. 1.

a significantly improved protocol; this is achieved by the high-level discrete optimization which selects the order of the operators in the sequence. However, we can do better: since  $|\psi_i\rangle$  is a product state while  $|\psi_*\rangle$  is not, and because  $H_3$  is a sum of single-particle terms, in order to create the target many-body correlations using a fast dynamical process, we also include the two-body first-order gauge potential terms  $H_4 = X|Y$  and  $H_5 = Y|Z$ : this results in a nonadiabatic evolution that prepares the interacting ground state to an excellent precision [Fig. 1 (red diamonds)].

In Ref. [42], it was shown that, in the integrable limit  $h_z = 0$ , one can prepare the ground state of the system at the critical point using a circuit of depth q = 2N with conventional QAOA. Albeit for the specific initial and target states chosen, we find that it only takes CD-QAOA a depth of q = 3 to reach the target ground state, independent of the system size N [43]. This result, though model-dependent, may come as a surprise at first sight, given that the mixed field Ising chain is a quantum chaotic system without a closed-form solution which makes it susceptible to heating away from the adiabatic limit.

Our data also reveals a finite many-body QSL at  $T_{\rm QSL} \approx 4.5$ . Importantly, this QSL appears insensitive to the system size to a very good approximation [Fig. 2], and we expect it to persist in the thermodynamic limit. The absence of a finite QSL in conventional QAOA in the mixed-field Ising chain suggests that the observation of a QSL using CD-QAOA depends on the specific set of unitaries related to the variational gauge potential, showcasing the utility of our ansatz for many-body control. Remarkably, we find an almost perfect system-size collapse of the target state energy density curves as a

function of the total protocol duration T. In Sec. VI, we explore this feature and demonstrate the ability of the RL agent to learn on small system sizes and subsequently generalize its knowledge to control bigger systems with exponentially larger Hilbert spaces.

CD-QAOA performs successfully on the nonintegrable spin-1/2 mixed-field Ising chain, for a circuit depth as short as q=3. This shows an advantage of our ansatz, when compared to conventional QAOA. However, the small size of the sequence space,  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1}=80$  at  $|\mathcal{A}|=5$ , poses a natural question regarding the necessity of using sophisticated search algorithms, such as RL, to find the control sequence. We now show that this is a peculiarity of the physical system, as we turn our attention to a larger sequence space.

#### B. Heisenberg Spin-1 Chain

The eight-dimensional spin-1 group SU(3) provides a significantly larger space of gauge potential terms to build the optimal protocol from. We consider a total of  $|\mathcal{A}| = 9$  unitaries: five are generated by the imaginary-valued terms in the gauge potential series: Y, XY, YZ, X|Y, Y|Z [cf. Table. I], plus the two real-valued QAOA operators  $H_1$  and  $H_2$ , which build the Hamiltonian  $H = H_1 + H_2$  whose ground state we target [Eq. (4)], and the two real-valued Hamiltonian terms X|X and Z. At q=18, this amounts to  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1} \approx 10^{16}$  possible sequences. The exponential scaling of the sequence space size with q renders applying exhaustive search algorithms infeasible, and justifies the use of sophisticated algorithms, such as RL.

The (anisotropic) spin-1 Heisenberg model reads as:

$$H = H_1 + H_2,$$

$$H_1 = J \sum_{j=1}^{N} (S_{j+1}^x S_j^x + S_{j+1}^y S_j^y), \quad H_2 = \Delta \sum_{j=1}^{N} S_{j+1}^z S_j^z,$$
(4)

with the spin exchange coupling J=1 set as energy unit, and  $\Delta$  – the anisotropy parameter; we use periodic boundary conditions and work in the ground state sector of zero momentum and positive parity, defined by the projector  $\mathcal{P}$ . In the thermodynamic limit, this model features a rich ground state phase diagram including ferromagnetic (FM,  $\Delta/J \ll -1$ ), XY  $(-1 \lesssim \Delta/J \lesssim 0)$ , topological/Haldane (0  $\lesssim \Delta/J \lesssim 1$ ), and antiferromagnetic (AFM,  $\Delta/J \gg 1$ ) order [44], with phase transitions belonging to different universality classes [45–47]. While the FM, XY, and AFM states are characterized by a local order parameter, the gapped Haldane state has topological order not captured by Landau-Ginzburg theory. We consider the AFM initial state  $|\psi_i\rangle = \mathcal{P} |\uparrow\downarrow\uparrow\downarrow\cdots\rangle$ , and target the ground states of Eq. (4) deep in the FM, XY, and Haldane phases, where system-size effects are the smallest. Because CD-QAOA is not restricted to adiabatic evolution, the conventional paradigm of a closing spectral gap when transferring the population between

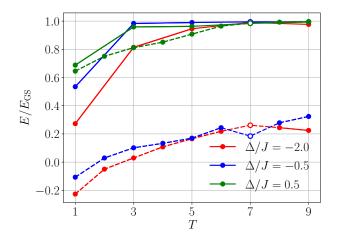


FIG. 3. Heisenberg spin-1 chain: energy minimization against protocol duration T using conventional QAOA (dashed lines) and CD-QAOA (solid lines) for three different states. We start from the AFM state  $|\psi_i\rangle = \mathcal{P} |\uparrow\downarrow\uparrow\downarrow \cdots\rangle$  and target three different parameter regimes, corresponding to the FM  $(\Delta/J = -2.0)$  state, XY  $(\Delta/J = -0.5)$ , and Haldane  $(\Delta/J=0.5)$  states, respectively. CD-QAOA outperforms conventional QAOA (p=q/2), more notably in the FM and XY targets where it allows us to reach close to the target state using a short protocol duration. The empty symbols mark the duration at which we show the evolution of the system in Fig. 20. The alternating unitaries for conventional QAOA are generated by  $\mathcal{A}_{QAOA} = \{H_1, H_2\}$  [cf. Eq. (4)]; for CD-QAOA, we extend this set using adiabatic gauge potential terms to  $\mathcal{A}_{\text{CD-QAOA}} = \{H_1, H_2, Z, X | X; Y, XY, YZ, X | Y - XY, Y | Z - XY,$ YZ. The circuit depths are q=28 ( $\Delta/J=-2.0$ ), q=18 $(\Delta/J = -0.5)$  and q = 18  $(\Delta/J = 0.5)$ . The cardinality of the CD-QAOA sequence space is  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1} \approx 10^{16}$  at q=18. The system size is N=8, where  $\dim(\mathcal{H})=498$ .

two states displaying different order, does not apply in our non-equilibrium setup, even in the thermodynamic limit.

Figure 3 shows a comparison between conventional QAOA with alternating sequence between the Hamiltonians  $H_1$  and  $H_2$ , and CD-QAOA. We find that CD-QAOA shows superior performance for all three ordered ground states: while the gain over conventional QAOA for the Haldane state is already a faster protocol, we clearly see how the gauge potential terms can prove essential for reaching the ground state in the FM and XY phases within the available durations. Note that the FM target state is doubly degenerate, and minimizing the energy, it ends up in an arbitrary superposition within the ground state manifold. Interestingly, we do not identify any distinction from preparing states with long-range and topological order, presumably due to the small system sizes that we reach in our classical simulation.

The CD-QAOA protocol sequences found by the RL agent have peculiar structures [App. F2]: some of them resemble closely the alternating sequence of conventional QAOA, with the notable difference of applying additional

unitaries to rotate the state to a suitable basis, either at the beginning or at the end of the sequence. While this is formally equivalent to starting from or targeting a rotated state, the rotations use two-body operators; hence, the resulting basis does not coincide with any of the distinguished  $S^x$ ,  $S^y$  and  $S^z$  directions. Variationally determining such effective bases demonstrates vet another advantage offered by the CD-QAOA ansatz. Another kind of encountered sequence contains two different sets of alternating unitaries, similar to two independent QAOA ansatzes concatenated one after the other. Finally, for those values of T, where CD-QAOA and QAOA have the same performance, we have also observed that CD-QAOA finds precisely the QAOA sequence. In this case, conventional QAOA already generates the shortest path, and the extra gauge potential terms to second-order do not give any advantage; a better performance might be expected when the three- and four-body higher-order terms from the gauge potential series are included.

Similar to other optimal control algorithms, RL agents typically find local minima of the optimization landscape; thus, there is no guarantee that the CD-QAOA protocols provide global optimal solutions; however, these sequences can serve as an inspiration to build future variational ansatzes tailored for many-body systems.

### C. Lipkin-Meshkov-Glick Model

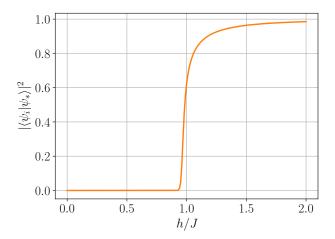
The non-integrable character of the previously discussed models precludes us from applying CD-QAOA with a large number of degrees of freedom, since reliably simulating their dynamics on a classical computer is prohibitively expensive. In order to study the behavior of CD-QAOA in a large enough system which also features a quantum phase transition, we now turn our attention to an exactly solvable many-body system.

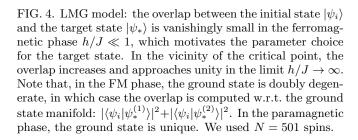
The Lipkin-Meshkov-Glick (LMG) Hamiltonian [48] describes spin-1/2 particles on a fully-connected graph of N sites:

$$H = H_1 + hH_2,$$

$$H_1 = -\frac{J}{N} \sum_{i,j=1}^{N} S_i^x S_j^x, \quad H_2 = \sum_{j=1}^{N} \left( S_j^z + \frac{1}{2} \right), \quad (5)$$

where J is the uniform interaction strength and h the external magnetic field. In the thermodynamic limit,  $N \to \infty$ , the system undergoes a quantum phase transition at  $h_c/J=1$  between a ferromagnetic (FM) ground state in the x-direction for  $h/J \ll 1$ , and a paramagnetic ground state for  $h/J \gg 1$ . The spectral gap  $\Delta_{\rm LMG}$  between the ground state and excited states closes as  $\Delta_{\rm LMG}(h_c) \sim N^{-1/3}$  at the critical point [49]. Realizing the LMG model is within the scope of present-day experiments with ultracold atoms [50, 51]; therefore, developing fast ground state preparation techniques can prove useful in practice.





Defining the total spin operators as  $S^{\alpha} = \sum_{j=1}^{N} S_{j}^{\alpha}$ , the Hamiltonian takes the form  $H = -J/N (S^{x})^{2} + h (S^{z} + N/2)$ . Hence, the total spin is conserved,  $[H, \mathbf{S} \cdot \mathbf{S}] = 0$ , and the ground state symmetry sector contains a total of N+1 states, i.e.  $\dim(\mathcal{H}) = N+1$ , which allows us to simulate large system sizes.

Our goal is, starting from the z-polarized paramagnetic initial state,  $|\psi_i\rangle=|\downarrow\downarrow\cdots\rangle$ , to target an arbitrary superposition in the doubly-degenerate FM ground state manifold, at fixed values of the external field h/J which controls the magnitude of the transversal fluctuations on top of the ferromagnetic order. Figure 4 shows that the overlap of the initial and target states is vanishingly small in the FM phase, and approaches quickly unity across the critical point into the paramagnetic phase. Therefore, we choose to prepare ground states in the FM phase where the problem naturally appears more difficult.

Figure 5 shows a comparison between CD-QAOA and QAOA on the LMG model at h/J=0.5 for N=501 spins [more h/J values are shown in App. F3]. First, note the superior performance of CD-QAOA, as compared to conventional QAOA in a range or short durations T in the nonadiabatic driving regime. We applied CD-QAOA with two different sets of generators:  $\mathcal{A} = \{H_1, H_2; Y\}$ , and  $\mathcal{A}' = \{H_1, H_2; Y, \hat{X}Y, \hat{Z}Y\}$  [cf. Table I] and found that, for the LMG model, the higher-order two-body terms  $\hat{X}Y, \hat{Z}Y$  do not offer any advantage deep in the FM phase. This observation can be understood as follows: to turn the z-polarized initial state into the x-ferromagnet, it is sufficient to perform a rotation about the y-axis,

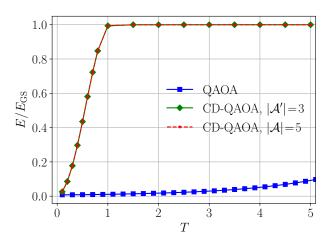


FIG. 5. LMG model: energy minimization against protocol duration T using conventional QAOA (blue square) and CD-QAOA (red dashed line, green solid line). We start from the z-polarized state  $|\psi_i\rangle = |\downarrow\downarrow \cdots\rangle$  and target ground state of LMG Hamiltonian (5). CD-QAOA significantly outperforms conventional QAOA for short durations. The alternating unitaries for conventional QAOA are generated by  $\mathcal{A}_{\text{QAOA}} = \{H_1, H_2\}$  [cf. Eq. (5)]; for CD-QAOA, we extend this set using adiabatic gauge potential terms to  $\mathcal{A}_{\text{CD-QAOA}} = \{H_1, H_2; Y, \hat{XY}, \hat{YZ}\}$  and  $\mathcal{A'}_{\text{CD-QAOA}} = \{H_1, H_2; Y\}$ . The external field is h/J = 0.5; the circuit depth is q = 8, and the system size is N = 501, where effective Hilbert dimension  $\dim(\mathcal{H}) = 502$ .

which coincides precisely with the single-body term in the gauge potential series expansion [cf. App. E1c]. Indeed, for all protocol durations smaller than the quantum speed limit,  $T < T_{\rm QSL}$ , the RL agent finds that the optimal protocol consists of a single Y-rotation, while for  $T \geq T_{\rm QSL}$  the optimal protocol is degenerate, and typically involves the various terms from  $\mathcal{A}$ . This finding allows us to extract the QSL as a function of the external field h, cf. Fig. 6.

Close to the critical point  $h_c$ , we observe strong sensitivity in the best found protocols to system-size effects, and a single Y-rotation is no longer optimal below the QSL. Interestingly, at the critical point (and in the paramagnetic phase), the optimal protocol is given by QAOA: in this regime, despite the larger set of terms  $\mathcal{A}$  we use in CD-QAOA, the RL agent correctly identifies the sequence of alternating  $H_1$  and  $H_2$  terms as optimal, which shows the versatility of CD-QAOA: the algorithm can always select a smaller effective subspace of actions when this is advantageous in the parameter regime of interest.

# V. COMPARISON WITH COUNTER-DIABATIC DRIVING

To compare and contrast the CD-QAOA ansatz with CD and adiabatic driving [34], consider the driven spin-1

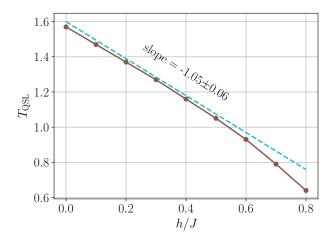


FIG. 6. LMG model: Quantum speed limit,  $T_{\rm QSL}$ , as a function of the transverse field h, for a target state in the ferromagnetic phase. At h/J=0, we have  $T_{\rm QSL}=\pi/2$ , which is the angle required to turn the z-polarized initial state into the x-ferromagnet. For finite h/J quantum fluctuations in the target ferromagnetic ground state decrease the angle required to transfer the population from the initial state, which results in a smaller value of  $T_{\rm QSL}$ . The dashed cyan line is a least squares fit for small values of h/J, suggesting the behavior  $T_{\rm QSL}(h)=-h/J+\pi/2+\mathcal{O}(h^2)$ . We used N=501 spins.

Ising model [52]:

$$H(\lambda) = \lambda(t)H_1 + H_2,$$

$$H_1 = \sum_{j=1}^{N} JS_{j+1}^z S_j^z + h_x S_j^x, \qquad H_2 = \sum_{j=1}^{N} h_z S_j^z,$$
(6)

where  $\lambda(t) = \sin^2\left(\frac{\pi t}{2T}\right)$ ,  $t \in [0,T]$ , is a smooth protocol satisfying the boundary conditions for CD driving:  $\lambda(0) = 0$ ,  $\lambda(T) = 1$ ,  $\dot{\lambda}(0) = 0 = \dot{\lambda}(T)$ . The initial state is the ground state at t = 0, i.e.  $|\psi_i\rangle = |\downarrow \cdots \downarrow\rangle$ , while the target state is the ground state of the Ising model at t = T for  $h_z/J = 0.809$  and  $h_x/J = 0.9045$ . Unlike the setup in Sec. IV A, adiabatic state preparation following the protocol  $\lambda(t)$ , suggests using the QAOA generators  $\mathcal{A}_{\text{QAOA}} = \{H_1, H_2\}$ .

Figure 7 shows a comparison between different methods using the best found energy density (main figure), and the corresponding many-body fidelity (inset). Let us focus on CD-QAOA and QAOA first. As expected, CD-QAOA (red) performs better for short durations T, since it contains conventional QAOA (red) as an ansatz, i.e.  $\mathcal{A}_{\text{QAOA}} \subsetneq \mathcal{A}_{\text{CD-QAOA}}$ . We emphasize that such a performance is not guaranteed in practice, since it is conceivable that the RL agent gets stuck in a local minimum associated with lower energy than the QAOA solution [App. D], e.g., if the deep autoregressive network architecture is not expressive enough, or if the learning rate schedules are not well-tuned to the problem. Unlike the spin-1/2 Ising model, here we cannot clearly identify a finite QSL, as the CD-QAOA energy keeps improving with

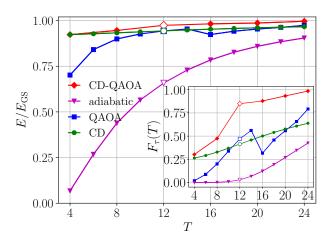


FIG. 7. Spin-1 Ising model: energy minimization and the corresponding many-body fidelity [inset] against different protocol duration T for four different optimization methods: CD-QAOA (red line), conventional QAOA (blue line), variational gauge potential (green) and adiabatic evolution (magenta). The empty symbols mark the duration for which the evolution of physical quantities is shown in Fig. 24. The initial and target states are  $|\psi_i\rangle = |\downarrow \cdots \downarrow\rangle$  and  $|\psi_*\rangle = |\psi_{\rm GS}(H)\rangle$  for  $h_z/J = 0.809$  and  $h_x/J = 0.9045$ . The alternating unitaries for conventional QAOA are generated by  $\mathcal{A}_{QAOA} = \{H_1, H_2\}$  [cf. Eq. (6)]; for CD-QAOA, we extend this set using adiabatic gauge potential terms to  $\mathcal{A}_{\text{CD-QAOA}} = \{H_1, H_2; Y, XY, YZ, X|Y, Y|Z\}.$  The variational gauge potential in CD driving uses all five imaginaryvalued gauge potentials  $\{Y, XY, YZ, X|Y, Y|Z\}$ . The CDand adiabatic driving simulations are both based on the smooth protocol function  $\lambda(t) = \sin^2\left(\frac{\pi t}{2T}\right)$ , with a timediscretization step  $\Delta t = 0.2$ . The value of q = 20 and the size of sequence space is  $|\mathcal{A}|(|\mathcal{A}|-1)^{q-1} \approx 10^{15}$ . The system size is N=8, where dim $(\mathcal{H})=498$ .

increasing circuit depth q [App. A].

To construct the counter-diabatic Hamiltonian  $H_{\rm CD} \approx$  $H(\lambda) + \lambda \mathcal{X}(\{\beta_i\})$  for Eq. (6), we make a variational ansatz [34] for the gauge potential  $\mathcal{X}$ , and solve for the optimal parameters  $\beta_i$  numerically [App. E]. We note the following differences between this approach and CD-QAOA: (i) the variational gauge potential depends on time t continuously, which requires further discretization when performing a gate-based implementation. (ii) the number of variational parameters in the standard variational gauge potential method is  $N_T|\mathcal{A}|$  with  $N_T$  the number of steps used to discretize the time interval [0, T]; instead, in CD-QAOA, we have q variational parameters. (iii) the variational gauge potential method does not constrain the magnitude of the variational coefficients  $\beta_i$ , and hence the time-averaged norm of  $H_{\rm CD}$  over the protocol can grow indefinitely; especially for short durations T this typically gives a higher fidelity. By contrast, in CD-QAOA the time-averaged norm of the unitary generators  $\alpha_j H_j$  summed along the sequence, is kept bounded via the constraint  $\sum_j \alpha_j = T$ . Nonetheless, in practice,

T	$E/E_{\rm GS}~[N=10]$				
	CD-QAOA	CD	QAOA	adiabatic	
4	0.943837	0.923199	0.79534	0.067807	
8	0.961383	0.933067	0.93386	0.438856	
12	0.990415	0.942857	0.96275	0.658182	

TABLE II. Spin-1 Ising model: comparison of the best obtained energy ratio  $E/E_{\rm GS}$  after optimization, for four different optimization methods: CD-QAOA, variational CD driving, conventional QAOA, and adiabatic evolution, at T=4,8,12 for N=10 qutrits, where  $\dim(\mathcal{H})=3219$ . The remaining setup and parameters are the same as in Fig. 7.

we find that these norms are on the same order of magnitude for all methods considered [App. F 4].

As anticipated, Fig. 7 shows that CD driving performs better than adiabatic driving, and the two agree in the limit of large T. Moreover, we see explicitly that the CD and QAOA solutions are far from the adiabatic regime. Not surprisingly, CD driving outperforms conventional QAOA for small T, as it can increase the values of the variational parameters (and with it the norm) indefinitely. However, CD-QAOA consistently outperforms CD driving in the entire T-range; the contrast is especially pronounced in the many-body fidelity [Fig. 7, inset]. CD-QAOA makes use of the variational power of QAOA, combining it with physics-motivated input from CD driving.

Table II shows a comparison with the best obtained energies for N=10 spin-1 particles (qutrits): the superior performance of CD-QAOA remains despite the exponentially growing Hilbert space size. Reaching significantly larger system sizes is infeasible with the present-day computational power: we note that this a feature of the quantum system rather than a drawback of CD-QAOA, cf. discussion on LMG model in Sec. IV C.

We emphasize that CD-QAOA features some important advantages as compared to CD driving: (1) Due to the nested commutators in the definition of time-ordered exponentials, the QAOA dynamics can effectively implement total unitaries  $U(\{\alpha_j\}_{j=1}^q, \tau)$  generated by effective non-local operators; therefore, CD-QAOA can, in principle, realize a nonlocal effective Hamiltonian as an approximation to the true CD Hamiltonian, thereby overcoming convergence issues related to operator-valued series expansions. (2) CD-QAOA lifts the boundary constraint present in adiabatic and CD driving where the initial and target Hamiltonians are eigenstates of H(0) and H(1), respectively; an interesting open question is whether a local effective Hamiltonian exists, which captures the evolution of the system in this case. Examining the evolution of the entanglement entropy and other local observables induced by the optimal protocol, suggests that this is indeed the case [App. F4]. (3) One can add any control unitary to the set A, not just terms related to gauge potentials: CD-QAOA has high flexibility to accommodate

experimental constraints. (4) determining the variational gauge potential in CD-driving requires using the exact ground state in order to minimize the action, which can be a significant drawback when the ground state is not known or cannot be computed.

# VI. TRANSFER LEARNING AND GENERALIZATION OF THE RL ALGORITHM TO DIFFERENT SYSTEM SIZES

The scale collapse in the energy density of the spin-1/2 Ising model presents a testbed for the transfer learning capabilities of RL. In transfer learning, the RL agent learns to control one physical system, and is then used to manipulate another. In our case, the two systems are given by the same Ising model at two different system sizes. Note that transfer learning would have not been possible, had we defined the learning problem using the full quantum states, because the latter are vectors in Hilbert space whose size grows exponentially with N.

To apply transfer learning, consider first a fixed protocol duration T. For every fixed system size N, we first train a different RL agent. Next, we build the set of protocols across all system sizes, found by these agents, and determine the number of unique protocols [cf. legend in Fig. 8. Finally, we apply all unique protocols to all system sizes available, and store the energy densities they result in. This leaves us with a set of energy density values for every fixed T. The error bars in Fig. 8 show the best and the worst protocols over this set. Observe that, below the QSL, there are only a few points T where the best control protocol is the same across all system sizes. Transfer learning works well, as can be seen by the small error bars. In this regime, the RL agent generalizes its knowledge and learns universal features of the protocol, required to control the Ising model. In contrast, for  $T > T_{\rm QSL}$ , there are many more protocols giving approximately similar ground state energies. While the corresponding energies are similar in value, the agent does not generalize. Nevertheless, we checked that, in this regime, training on smaller system sizes still provides a useful pre-training procedure for learning on larger systems.

# VII. DISCUSSION/OUTLOOK

We analyzed many-body ground state preparation using unitary evolution in the spin-1/2 Ising model, the spin-1 anisotropic Heisenberg and Ising models, and the fully connected LMG spin-1/2 model. We introduced the CD-QAOA ansatz: an RL agent optimizes the order of unitaries in the protocol sequence, generated from terms in the adiabatic gauge potential series, and obtains short high-fidelity protocols away from the adiabatic regime. The resulting algorithm combines the strength of continuous and discrete optimization into a unified and versatile control framework. We find that our

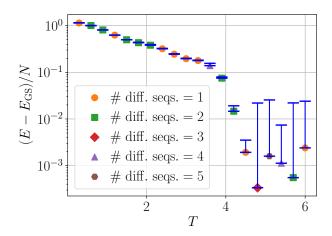


FIG. 8. Spin-1/2 Ising model: Protocol generalization across various system sizes. The marker types show the number of different protocols found by the RL agent at a fixed T across all system sizes N=6,10,12,14,16 and 18. Each protocol is applied to every system size N at a fixed T which results in a set of cost function values; the error bars designate the range between the largest and smallest cost function value. The parameters are the same as in Fig. 1.

CD-QAOA ansatz outperforms consistently both conventional QAOA, and variational CD driving across different models and protocol durations. An interesting open question is whether one can use CD-QAOA to find a nonlocal approximation to the variational gauge potential itself, which is beyond the scope of asymptotic series expansions. Another straightforward application of CD-QAOA would be imaginary time evolution [53].

For the nonintegrable spin-1/2 Ising chain, we reveal the existence of a finite quantum speed limit. Moreover, we find a remarkable system-size collapse of the energy curves suggesting that the sequences found by the agent hold in the thermodynamic limit; this is corroborated by numerical experiments on transfer learning which demonstrate that one can train the agent on one system size while it generalizes to larger systems. In the Heisenberg spin-1 system, CD-QAOA allows preparing long-range and topologically ordered ground states, even when the initial state does not belong to the phase of the target state. The optimal protocols found by the RL agent contain nontrivial basis rotations, intertwined with alternating QAOA-like subsequences, suggesting new ansätze for more efficient variants of CD-QAOA. Numerical studies of nonequilibrium quantum many-body systems, in turn, suffer from limitations related to the exponentially large

dimension of the underlying Hilbert space: future work can investigate dynamics beyond exact evolution.

Compared to conventional QAOA, using terms from the variational gauge potential series has higher expressivity, which results in much shorter, yet better performing, circuits. This method can be used, e.g., to reduce the cumulative error in quantum computing devices. However, gauge potential terms are not always easy to realize in experiments since they implement imaginary-valued terms which break time-reversal symmetry; that said, it is often possible to generate such terms using auxiliary real-valued operators via a generalization of the Euler angles, or by means of change-of-frame transformations [34]. Moreover, as we have demonstrated, CD-QAOA admits non-gauge potential terms as building blocks for control sequences, e.g., universal gate sets. Other experimental constraints, such as the presence of drift terms, which cannot be switched off, can also be conveniently incorporated by redefining the set of unitaries A.

Finally, let us remark that RL provides only one possible set of algorithms to explore the exponentially large space of protocol sequences; in practice, one can apply other discrete optimization techniques, e.g. genetic algorithms and search algorithms like Monte-Carlo Tree Search (MCTS).

Acknowledgments.—We wish to thank A. Polkovnikov, Dong An and Yulong Dong for valuable discussions. This work was partially supported by the Department of Energy under Grant No. DE-AC02-05CH11231, No. DE-SC0017867 and by a Google Quantum Research Award (L.L., J.Y.), and by the NSF Quantum Leap Challenge Institute (QLCI) program through grant number OMA-2016245 (L.L.). M.B. was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, under the Accelerated Research in Quantum Computing (ARQC) program, the Quantum Algorithm Teams Program, the U.S. Department of Energy under cooperative research agreement DE-SC0009919, the Emergent Phenomena in Quantum Systems initiative of the Gordon and Betty Moore Foundation, and the Bulgarian National Science Fund within National Science Program VIHREN, contract number KP-06-DV-5. We used SLSQP implemented in SciPy for the QAOA solver, NumPy, and TensorFlow and Tensor-Flow Probability for the deep learning simulations; we used Quspin for simulating the dynamics of the quantum systems [54, 55]. The authors are pleased to acknowledge that the computational work reported on in this paper was performed on Savio3 Condo of Berkeley Research Computing (BRC).

M. Lewenstein, A. Sanpera, V. Ahufinger, B. Damski, A. Sen, and U. Sen, Ultracold atomic gases in optical lattices: mimicking condensed matter physics and beyond, Advances In Physics 56, 243 (2007).

<sup>[2]</sup> I. Bloch, J. Dalibard, and W. Zwerger, Many-body

physics with ultracold gases, Rev. Mod. Phys.  $\bf 80$ , 885 (2008).

<sup>[3]</sup> H. Häffner, C. F. Roos, and R. Blatt, Quantum computing with trapped ions, Physics reports **469**, 155 (2008).

<sup>[4]</sup> R. Blatt and C. F. Roos, Quantum simulations with

- trapped ions, Nature Physics 8, 277 (2012).
- [5] C. Monroe and J. Kim, Scaling the ion trap quantum processor, Science 339, 1164 (2013).
- [6] M. H. Devoret and R. J. Schoelkopf, Superconducting circuits for quantum information: an outlook, Science 339, 1169 (2013).
- [7] M. W. Doherty, N. B. Manson, P. Delaney, F. Jelezko, J. Wrachtrup, and L. C. Hollenberg, The nitrogenvacancy colour centre in diamond, Physics Reports 528, 1 (2013).
- [8] R. Schirhagl, K. Chang, M. Loretz, and C. L. Degen, Nitrogen-vacancy centers in diamond: nanoscale sensors for physics and biology, Annual review of physical chemistry 65, 83 (2014).
- [9] F. Casola, T. van der Sar, and A. Yacoby, Probing condensed matter physics with magnetometry based on nitrogen-vacancy centres in diamond, Nature Reviews Materials 3, 17088 (2018).
- [10] Z.-L. Xiang, S. Ashhab, J. Q. You, and F. Nori, Hybrid quantum circuits: Superconducting circuits interacting with other quantum systems, Rev. Mod. Phys. 85, 623 (2013).
- [11] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. Brandao, D. A. Buell, et al., Quantum supremacy using a programmable superconducting processor, Nature 574, 505 (2019).
- [12] N. Khaneja, T. Reiss, C. Kehlet, T. Schulte-Herbrüggen, and S. J. Glaser, Optimal control of coupled spin dynamics: design of NMR pulse sequences by gradient ascent algorithms, Journal of Magnetic Resonance 172, 296 (2005).
- [13] T. Caneva, T. Calarco, and S. Montangero, Chopped random-basis quantum optimization, Phys. Rev. A 84, 022326 (2011).
- [14] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O'brien, A variational eigenvalue solver on a photonic quantum processor, Nature communications 5, 4213 (2014).
- [15] L. Zhou, S.-T. Wang, S. Choi, H. Pichler, and M. D. Lukin, Quantum approximate optimization algorithm: Performance, mechanism, and implementation on near-term devices, Physical Review X 10, 021067 (2020).
- [16] M. Demirplak and S. A. Rice, Adiabatic population transfer with control fields, The Journal of Physical Chemistry A 107, 9937 (2003).
- [17] M. Berry, Transitionless quantum driving, Journal of Physics A: Mathematical and Theoretical 42, 365303 (2009).
- [18] M. Kolodrubetz, D. Sels, P. Mehta, and A. Polkovnikov, Geometry and non-adiabatic response in quantum and classical systems, Physics Reports 697, 1 (2017).
- [19] M. Bukov, D. Sels, and A. Polkovnikov, Geometric speed limit of accessible many-body state preparation, Phys. Rev. X 9, 011034 (2019).
- [20] We focus on pure states, although the cost function can trivially be generalized to mixed states.
- [21] V. Jurdjevic and H. J. Sussmann, Control systems on lie groups, Journal of Differential equations 12, 313 (1972).
- [22] L. Zhu, H. L. Tang, G. S. Barron, F. A. Calderon-Vargas, N. J. Mayhall, E. Barnes, and S. E. Economou, An adaptive quantum approximate optimization algorithm for solving combinatorial problems on a quantum computer, arXiv preprint arXiv:2005.10258v2 (2020).
- [23] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Bab-

- bush, and H. Neven, Barren plateaus in quantum neural network training landscapes, Nature Communications 9, 4812 (2018).
- [24] M. Cerezo, A. Sone, T. Volkoff, L. Cincio, and P. J. Coles, Cost function dependent barren plateaus in shallow parametrized quantum circuits, Nature Communications 12, 1791 (2021).
- [25] E. Grant, L. Wossnig, M. Ostaszewski, and M. Benedetti, An initialization strategy for addressing barren plateaus in parametrized quantum circuits, Quantum 3, 214 (2019).
- [26] P. Huembeli and A. Dauphin, Characterizing the loss landscape of variational quantum circuits, Quantum Science and Technology 6, 025011 (2021).
- [27] Considering  $\tau_j$  as choice of unitaries, we impose the extra constraint that, even though unitaries can be repeated in the sequence  $\tau$ , the same unitary cannot appear consecutively (or else one can combine the two corresponding choices  $\tau_i$  into a single variable).
- [28] A. G. R. Day, M. Bukov, P. Weinberg, P. Mehta, and D. Sels, Glassy phase of optimal quantum control, Phys. Rev. Lett. 122, 020601 (2019).
- [29] M. Bukov, A. G. R. Day, P. Weinberg, A. Polkovnikov, P. Mehta, and D. Sels, Broken symmetry in a two-qubit quantum control landscape, Phys. Rev. A 97, 052114 (2018).
- [30] A similar procedure appeared recently in Ref. [?], although they considered a different problem setup with greedy or beam search algorithm.
- [31] In principle, one can use any optimizer which allows constraining the sum  $\sum_{j} \alpha_{j} = T$ .
- [32] N. Lacroix, C. Hellings, C. K. Andersen, A. D. Paolo, A. Remm, S. Lazar, S. Krinner, G. J. Norris, M. Gabureac, J. Heinsoo, A. Blais, C. Eichler, and A. Wallraff, Improving the performance of deep quantum optimization algorithms with continuous gate sets, PRX Quantum 1, 110304 (2020).
- [33] Y. Ding, Y. Ban, J. D. Martín-Guerrero, E. Solano, J. Casanova, and X. Chen, Breaking adiabatic quantum control with deep learning, Physical Review A 103, L040401 (2021).
- [34] D. Sels and A. Polkovnikov, Minimizing irreversible losses in quantum systems by local counterdiabatic driving, Proceedings of the National Academy of Sciences 114, E3909 (2017).
- [35] A. Hartmann and W. Lechner, Rapid counter-diabatic sweeps in lattice gauge adiabatic quantum computing, New Journal of Physics 21, 043025 (2019).
- [36] J. Wurtz, P. W. Claeys, and A. Polkovnikov, Variational schrieffer-wolff transformations for quantum many-body dynamics, Phys. Rev. B 101, 014302 (2020).
- [37] N. N. Hegade, K. Paul, Y. Ding, M. Sanz, F. Albarrán-Arriagada, E. Solano, and X. Chen, Shortcuts to adiabaticity in digitized adiabatic quantum computing, Physical Review Applied 15 (2021).
- [38] M. Pandey, P. W. Claeys, D. K. Campbell, A. Polkovnikov, and D. Sels, Adiabatic eigenstate deformations as a sensitive probe for quantum chaos, Physical Review X 10, 041017 (2020).
- [39] Below, we sometimes abuse notation and set  $A = \{H_j\}$ , denoting the set of unitaries by their generators.
- [40] J. Wurtz and P. J. Love, Counterdiabaticity and the quantum approximate optimization algorithm, arXiv

- preprint arXiv:2106.15645 (2021).
- [41] G. Matos, S. Johri, and Z. Papić, Quantifying the efficiency of state preparation via quantum variational eigensolvers, PRX Quantum 2, 010309 (2021).
- [42] W. W. Ho and T. H. Hsieh, Efficient variational simulation of non-trivial quantum states, SciPost Phys. 6, 29 (2019).
- [43] The role of the RL algorithm is to decide which three out of the five unitaries  $U_j$  to apply and in which order.
- [44] We define 'order' in the context of phase transitions in condensed matter physics.
- [45] W. Chen, K. Hida, and B. C. Sanctuary, Ground-state phase diagram of s=1 XXZ chains with uniaxial single-ion-type anisotropy, Phys. Rev. B **67**, 104401 (2003).
- [46] F. Pollmann, A. M. Turner, E. Berg, and M. Oshikawa, Entanglement spectrum of a topological phase in one dimension, Phys. Rev. B 81, 064439 (2010).
- [47] A. Langari, F. Pollmann, and M. Siahatgar, Groundstate fidelity of the spin-1 heisenberg chain with single ion anisotropy: quantum renormalization group and exact diagonalization approaches, Journal of Physics: Condensed Matter 25, 406002 (2013).
- [48] H. Lipkin, N. Meshkov, and A. Glick, Validity of many-body approximation methods for a solvable model: (i). exact solutions and perturbation theory, Nuclear Physics 62, 188 (1965).
- [49] R. Botet and R. Jullien, Large-size critical behavior of infinitely coordinated systems, Phys. Rev. B 28, 3955 (1983).
- [50] H. Strobel, W. Muessel, D. Linnemann, T. Zibold, D. B. Hume, L. Pezzè, A. Smerzi, and M. K. Oberthaler, Fisher information and entanglement of non-gaussian spin states, Science 345, 424 (2014).
- [51] E. J. Davis, A. Periwal, E. S. Cooper, G. Bentsen, S. J. Evered, K. Van Kirk, and M. H. Schleier-Smith, Protecting spin coherence in a tunable heisenberg model, Phys. Rev. Lett. 125, 060402 (2020).
- [52] We deliberately use a different form in Eq. (6) as compared to Eq. (3); the former may appear more natural in quantum many-body physics, where the transverse-field Ising model  $H_1$  can be mapped to free fermions.
- [53] M. J. S. Beach, R. G. Melko, T. Grover, and T. H. Hsieh, Making trotters sprint: A variational imaginary time ansatz for quantum many-body systems, Phys. Rev. B 100, 094434 (2019).
- [54] P. Weinberg and M. Bukov, QuSpin: a Python Package for Dynamics and Exact Diagonalisation of Quantum Many Body Systems part I: spin chains, SciPost Phys. 2, 003 (2017).
- [55] P. Weinberg and M. Bukov, QuSpin: a python package for dynamics and exact diagonalisation of quantum many body systems. part II: bosons, fermions and higher spins, SciPost Physics 7, 020 (2019).
- [56] V. Dunjko and H. J. Briegel, Machine learning & artificial intelligence in the quantum domain: a review of recent progress, Reports on Progress in Physics 81, 074001 (2018).
- [57] P. Mehta, M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher, and D. J. Schwab, A high-bias, low-variance introduction to machine learning for physicists, Physics reports 810, 1 (2019).
- [58] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, Machine learning and the physical sciences, Rev. Mod. Phys. 91,

- 045002 (2019).
- [59] J. Carrasquilla, Machine learning for quantum matter, Advances in Physics: X 5, 1797528 (2020).
- [60] K. J. Sung, J. Yao, M. P. Harrigan, N. C. Rubin, Z. Jiang, L. Lin, R. Babbush, and J. R. McClean, Using models to improve optimizers for variational quantum algorithms, Quantum Science and Technology 5, 044008 (2020).
- [61] F. Schäfer, M. Kloc, C. Bruder, and N. Lörch, A differentiable programming method for quantum control, Machine Learning: Science and Technology 1, 035009 (2020).
- [62] F. Sauvage and F. Mintert, Optimal quantum control with poor statistics, PRX Quantum 1, 020322 (2020).
- [63] T. Fösel, S. Krastanov, F. Marquardt, and L. Jiang, Efficient cavity control with snap gates, arXiv preprint arXiv:2004.14256v1 (2020).
- [64] R.-B. Wu, X. Cao, P. Xie, and Y.-x. Liu, End-to-end quantum machine learning implemented with controlled quantum dynamics, Physical Review Applied 14, 064020 (2020).
- [65] F. Albarrán-Arriagada, J. C. Retamal, E. Solano, and L. Lamata, Measurement-based adaptation protocol with quantum reinforcement learning, Phys. Rev. A 98, 042315 (2018).
- [66] T. Fösel, P. Tighineanu, T. Weiss, and F. Marquardt, Reinforcement learning with neural networks for quantum feedback, Phys. Rev. X 8, 031084 (2018).
- [67] H. P. Nautrup, N. Delfosse, V. Dunjko, H. J. Briegel, and N. Friis, Optimizing quantum error correction codes with reinforcement learning, Quantum 3, 215 (2019).
- [68] R. S. Sutton and A. G. Barto, Reinforcement learning: An introduction (MIT press, 2018).
- [69] D. C. Rose, J. F. Mair, and J. P. Garrahan, A reinforcement learning approach to rare trajectory sampling, New Journal of Physics 23, 013013 (2021).
- [70] M. Y. Niu, S. Boixo, V. N. Smelyanskiy, and H. Neven, Universal quantum control through deep reinforcement learning, npj Quantum Information 5, 1 (2019).
- [71] M. August and J. M. Hernández-Lobato, Taking gradients through experiments: Lstms and memory proximal policy optimization for black-box quantum control, in *High Performance Computing* (Springer International Publishing, Cham, 2018) pp. 591–613.
- [72] R. Porotti, D. Tamascelli, M. Restelli, and E. Prati, Coherent transport of quantum states by deep reinforcement learning, Communications Physics 2, 61 (2019).
- [73] M. Bukov, Reinforcement learning for autonomous preparation of floquet-engineered states: Inverting the quantum kapitza oscillator, Physical Review B 98, 224305 (2018).
- [74] M. Bukov, A. G. R. Day, D. Sels, P. Weinberg, A. Polkovnikov, and P. Mehta, Reinforcement learning in different phases of quantum control, Phys. Rev. X 8, 031086 (2018).
- [75] M. Dalgaard, F. Motzoi, J. J. Sørensen, and J. Sherson, Global optimization of quantum dynamics with alphazero deep exploration, npj Quantum Information 6, 6 (2020).
- [76] J. Yao, M. Bukov, and L. Lin, Policy gradient based quantum approximate optimization algorithm, in *Mathematical and Scientific Machine Learning Conference* (MSML), 2020, Vol. 107 (PMLR, 2020) pp. 605–634.
- [77] M. M. Wauters, E. Panizon, G. B. Mbeng, and G. E. Santoro, Reinforcement learning assisted quantum opti-

- mization, Phys. Rev. Research 2, 033446 (2020).
- [78] S. Khairy, R. Shaydulin, L. Cincio, Y. Alexeev, and P. Balaprakash, Reinforcement-learning-based variational quantum circuits optimization for combinatorial problems, arXiv preprint arXiv:1911.04574v1 (2019).
- [79] A. Garcia-Saez and J. Riu, Quantum observables for continuous control of the quantum approximate optimization algorithm via reinforcement learning, arXiv preprint arXiv:1911.09682v1 (2019).
- [80] A. Bolens and M. Heyl, Reinforcement learning for digital quantum simulation, arXiv preprint arXiv:2006.16269v1 (2020).
- [81] It is also possible to define an RL framework for hybrid continuous-discrete control where optimization is entirely based on RL, cf. Ref. [90].
- [82] M. Bukov, Reinforcement learning for autonomous preparation of floquet-engineered states: Inverting the quantum kapitza oscillator, Phys. Rev. B 98, 224305 (2018).
- [83] D. Wu, L. Wang, and P. Zhang, Solving statistical mechanics using variational autoregressive networks, Physical Review Letters 122, 080602 (2019).
- [84] O. Sharir, Y. Levine, N. Wies, G. Carleo, and A. Shashua, Deep autoregressive models for the efficient variational simulation of many-body quantum systems, Phys. Rev. Lett. 124, 020503 (2020).
- [85] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal policy optimization algorithms, arXiv preprint arXiv:1707.06347v2 (2017).
- [86] B. D. Ziebart, Modeling purposeful adaptive behavior with the principle of maximum causal entropy (Carnegie Mellon University, USA, 2010).
- [87] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, arXiv preprint arXiv:1801.01290v2 (2018).
- [88] D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980v9 (2014).
- [89] J. Nocedal and S. Wright, *Numerical optimization* (Springer Science & Business Media, 2006).
- [90] J. Yao, P. Köttering, H. Gundlach, L. Lin, and M. Bukov, Noise-robust end-to-end quantum control using deep autoregressive policy networks, Mathematical and Scientific Machine Learning Conference (MSML), 2021, accepted, arXiv preprint arXiv:2012.06701v1 (2020).
- [91] T. Hatomura, Shortcuts to adiabaticity in the infiniterange ising model by mean-field counter-diabatic driving, Journal of the Physical Society of Japan 86, 094002 (2017).
- [92] Z. Mzaouali, R. Puebla, J. Goold, M. E. Baz, and S. Campbell, Work statistics and symmetry breaking in an excited-state quantum phase transition, Physical Review E 103, 032145 (2021).

# Appendix A: High-level optimization: Policy Gradient using Deep Autoregressive Networks

Recently, progress made in machine learning (ML) [56–59] has raised the question as to how we can harness such modern advances to improve techniques to manipulate quantum systems. Examples of ML applications include model-based optimization [60], differentiable programming [61] and Bayesian inference [62] quantum control, cavity control [63], designing quantum end-to-end learning schemes [64] and measurement-based adaptation protocols [65], as well as applications to quantum error-correction [66, 67].

Reinforcement learning (RL) algorithms [68, 69], such as policy gradient [70–72], Q-learning [73, 74] and AlphaZero [75], have recently attracted the attention of physicists, and in particular how they can be combined with physically motivated VQEs for improved performance. In RL, policy gradient has been proposed as an alternative optimizer for QAOA showcasing the robustness of RL-based optimization to both classical and quantum sources of noise [76]; a related study applied Proximal Policy Optimization (PPO) to prepare the ground state of the transverse-field Ising model [77]. The QAOA ansatz with policy gradient has been applied to efficiently find optimal variational parameters for unseen combinatorial problem instances on a quantum computer [78]; Qlearning was used to formulate QAOA into an RL framework to solve difficult combinatorial problems [79], and in the context of digital quantum simulation [80].

In the following, we introduce the details of the Reinforcement Learning algorithm used for the high-level optimization in this work.

# 1. Reinforcement Learning Basics

Reinforcement learning (RL) comprises a class of machine learning algorithms where an agent learns how to solve a given task through interactions with its environment using a trial-and-error approach [68]. It is based on a Markov Decision Process (MDP) defined by the tuple (S, A, p, R) where S and A represent the state and action spaces,  $p: \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$  defines the transition dynamics, and  $R: S \times A \rightarrow \mathbb{R}$  is the reward function that describe the environment. Let  $\pi(a_i|s_i)$ :  $\mathcal{A} \times \mathcal{S} \rightarrow [0,1]$  denote a stochastic policy that defines the probability distribution of choosing an action  $a_j \in \mathcal{A}$ given the state  $s_j \in \mathcal{S}$ . Rolling out the policy  $\pi(a_j|s_j)$ in the environment can also be viewed as sampling a trajectory  $\tau \sim \mathbb{P}^{\pi}(\cdot)$  from the MDP, where  $\mathbb{P}^{\pi}(\tau) =$  $p_0(s_1)\pi(a_1|s_1)p(s_2|s_1,a_1)\cdots\pi(a_q|s_q)p(s_{q+1}|s_q,a_q)$  is the probability for the trajectory  $\tau$  to occur, q sets the episode or trajectory length, and  $p_0$  is the distribution of the initial state; an example for a trajectory is  $\tau = (s_1, a_1, ..., a_q, s_{q+1})$ . The goal in RL is to find a

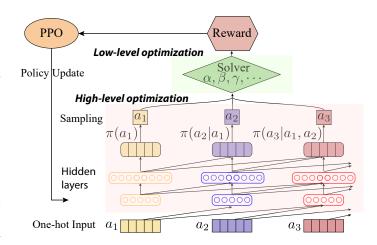


FIG. 9. Schematics of CD-QAOA with an autoregressive policy network. The ancestral sampling procedure used for training is displayed in Fig. 10. The details of the network structure and its training hyperparameters are shown in Table III.

policy that maximizes the expected return:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \mathbb{P}^{\pi}} \left[ \sum_{j=1}^{q} R(s_j, a_j) \right]. \tag{A1}$$

To maximize the expected return  $J(\theta)$ , we use policy gradient – an RL algorithm, which is (i) on-policy, i.e. trajectories have to be sampled from the current policy  $\pi_{\theta}$ :  $\pi = \pi_{\theta}$ , and (ii) model-free, i.e. the agent does not need to have a model for the environment dynamics: p(s'|s,a) is assumed unknown for the purpose of finding the optimal policy. Highly expressive function approximators, such as deep neural networks, help parametrize the policy using variational parameters  $\theta$ . Policy gradient gradually improves the expected return in a number of iterations (or training episodes), by increasing the probability for actions that lead to higher rewards, and decreasing the probability for actions that lead to lower rewards, until it reaches a (nearly) optimal policy.

We mention in passing that we use interchangeably the terms return and cost function (the latter being the negative of the former): the goal of the RL agent is thus to maximize the expected return, or to minimize the cost function.

# 2. Policy Gradient Reinforcement Learning for Quantum Many-Body Systems

**Actions:** To apply the reinforcement learning formalism to quantum control, we identify taking actions at each time step within a learning episode, with selecting unitaries one at a time within the circuit depth q. Choosing the same unitary at two consecutive time steps is prohibited because the same actions can be merged resulting in a lower effective circuit depth q-1. At the initial time

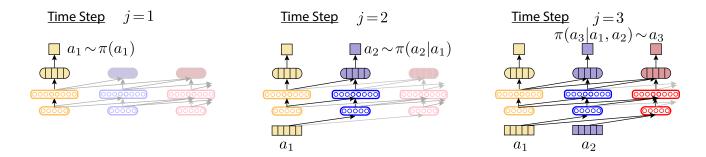


FIG. 10. The exact sampling algorithm for CD-QAOA with an autoregressive policy network, where faded nodes and connections represent unused nodes and connections. The action at each time step is generated sequentially, by computing its respective conditional categorical distribution, and sampling according to that. Notice that only a single column is processed at each time step, and in order to sample a complete sequence of actions in an episode one needs to make a forward pass through the network architecture q times.

step j=1, the quantum wavefunction is given by the initial state  $|\psi_i\rangle$ ; for each intermediate protocol step j, the action  $a_j=H_j$  is chosen according to the policy  $\pi_{\boldsymbol{\theta}}$ . Note that the RL agent only selects the generator  $H_j$  out of the set of available actions  $\mathcal{A}$  (or alternatively – which unitary to apply). In other words, unlike Ref. [76], the RL part of CD-QAOA is not concerned with finding the corresponding optimal duration  $\alpha_j$ ; one can think of this low-level continuous optimization as being part of the environment [cf. App. B][81]. At the end of the episode, the quantum state is evolved by applying the entire generated circuit  $U(\{\alpha_j\}_{j=1}^q, \tau)$  to the initial quantum state  $|\psi_i\rangle$ .

**States:** Since the initial state  $|\psi_i\rangle$  is fixed and thus the quantum state at any time step j is uniquely determined by the previous actions taken, here we represent the RL state by concatenating all the previous actions up to step j [82]. One reason for this is that, in many-body quantum systems, the number of components in the quantum state scales exponentially with the system size N, which quickly leads to a computational bottleneck for the simulation on classical computers. A second advantage of this choice is that the first layer of the underlying deep neural network architecture, which parametrizes the policy, will not depend on the system size N either, which allows the algorithm to handle a large number of degrees of freedom. Using the quantum state would not be viable on quantum computers either, because quantum states are unphysical mathematical constructs that cannot be measured. Therefore, we can simplify the form of trajectories to consist of actions only, e.g.  $\tau = (a_1, a_2, \dots, a_q)$ .

**Rewards:** The reward  $R_j = R(s_j, a_j)$  is chosen as the negative energy density at the end of the episode:

$$R_j = \left\{ \begin{array}{ll} 0, & \text{if } j < q \\ -E(\{\alpha_j\}_{j=1}^q, \tau)/N, & \text{if } j = q. \end{array} \right.$$

We use energy density, since it is an intensive quantity which has a well-defined limit when increasing the number of particles N. In all figures, we show the relative energy  $E/E_{\rm GS}$  for clarity (the ground state energy

 $E_{\rm GS}$  is typically negative in our models), but the RL agent is always trained with the (negative) energy density -E/N. Rewards can also be other observables or nonobservable quantities, such as the overlap squared between two quantum states (fidelity), or the entanglement entropy.

Notice that the reward is sparse: only at the end of the episode is the negative energy density given as a reward; there is no instantaneous reward during the sequence [and thus we can use interchangeably the terms reward and total return]. This is motivated by the quantum nature of the control problem, where a projective measurement results in a wavefunction collapse.

# 3. Policy Parametrization using an Autoregressive Neural Network

An essential part of the policy gradient algorithm is the definition of the policy  $\pi_{\theta}$ . It is common to parametrize the policy with a highly expressive function approximator, such as a neural network. In our setup, we use a deep autoregressive network, which has recently been used in physics applications of learning to generate samples from free energies in statistical mechanics models [83], and variational approximators for quantum many-body states [84]. This architecture is selected to incorporate causality by factorizing the total probability into a product of conditional probabilities:

$$\pi_{\theta}(a_1, a_2, \dots, a_q) = \pi_{\theta}(a_1) \prod_{j=2}^{q} \pi_{\theta}(a_j | a_1, \dots, a_{j-1}),$$
(A2)

where the marginal distribution  $\pi_{\theta}(a_1)$  and the conditional distribution  $\pi_{\theta}(a_j|a_1,\dots,a_{j-1})$  are discrete categorical distributions over  $\mathcal{A}$ . This kind of parametrization explicitly tells how the actions taken in the earlier steps of an episode affect the actions selected later on during the same episode. Such a causal requirement would not be necessary, had we used the full quantum

state, which would make the dynamics of the environment Markovian. Each of the conditional probabilities in Eq. (A2) can be modeled explicitly using the autoregressive neural network architecture, which naturally allows the policy to depend on all the previous actions only. The structure of the policy network is shown in Fig. 9, the sampling of the autoregressive policy is depicted in Fig. 10 and the hyperparameters of the algorithm (including the number of parameters) are given in Table III.

# 4. Training Procedure: Proximal Policy Optimization (PPO)

In each iteration of the policy gradient algorithm, a batch of sampled trajectories  $\{\tau^k\} = \{(a_1^k, \cdots, a_q^k)\}_{k=1}^M$  are rolled out (i.e. sampled) from the current policy, where M is the batch/sample size. Then, the return  $R(\tau^k)$  corresponding to trajectory  $\tau^k$  is computed as

$$R(\tau^k) = \sum_{j=1}^q R_j^k = -E(\{\alpha_j^k\}_{j=1}^q, \tau^k)/N.$$

To compute the energies, we use the low-level optimization to determine the best-estimated values of  $\alpha_j$ , given a sequence  $\tau$ , see App. B. To minimize the chance of getting stuck in a suboptimal local minimum, each sequence is evaluated multiple times, starting from a different initial realization for the  $\alpha_j$ -optimizer, and the best result is selected [App. D].

For every iteration, we can define three quantities for a fixed batch of samples: (i) mean reward (over the current batch), (ii) max reward (over the current batch), and (iii) history best (best-encountered reward over all the previous iterations). These quantities measure the performance of the learned policy, and are shown in Fig. 11. Figure 12 shows the scaling of these quantities for the spin-1 Ising chain, as a function of the episode length q. The performance of CD-QAOA increases because the action space for a larger value of q always contains as a subset the action space for a smaller q.

In order to improve the policy represented by the autoregressive network, the RL algorithm interacts with the quantum environment by querying the reward for samples from the current policy. Each trajectory is assigned a reward, once the simulation of the quantum dynamics is complete [note that, as of present date, the simulation may be more expensive if evaluated on a quantum computer]. Thus, it is advantageous to reduce the sample size needed to learn the policy, i.e., to improve the sample efficiency.

The vanilla policy gradient method is known for its poor data efficiency. Thus, we adopt Proximal Policy Optimization (PPO) [85], a more robust and sample-efficient policy gradient type algorithm. To be more specific, we

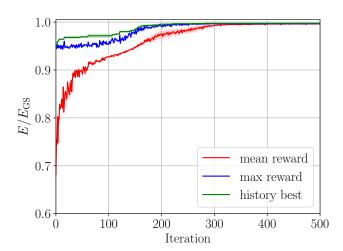


FIG. 11. Spin-1 Ising model: training curves for CD-QAOA with energy minimization as a cost function. The mean negative energy density (red) is computed for a sample generated using the policy at the current iteration; max (blue) is the maximum within the sample; the history best (green) is the best-encountered policy during the entire training process (i.e., considering all iterations). Each curve shows the average out of three simulations corresponding to three different seed values for the high level RL optimization; the fluctuations around the seed-averages are shown as a narrow shaded area. The total duration is T=28 and the number of spin-1 particles is N=8. The initial and target states are  $|\psi_i\rangle = |\downarrow \cdots \downarrow\rangle$  and  $|\psi_*\rangle = |\psi_{GS}(H)\rangle$  for  $h_z/J = 0.809$  and  $h_x/J = 0.9045$ . The CD-QAOA action space is  $\mathcal{A}_{\text{CD-QAOA}} = \{Z|Z+X,Z;Y,XY,YZ,X|Y,Y|Z\}$ , and we use q = 20.

use the following clipped objective function:

$$\mathcal{G}(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}_t}} \left[ \min \left\{ \rho_{\boldsymbol{\theta}}(\tau) A_{\boldsymbol{\theta}_t}(\tau), \right. \right.$$

$$\left. \text{clip} \left( \rho_{\boldsymbol{\theta}}(\tau), 1 - \epsilon, 1 + \epsilon \right) A_{\boldsymbol{\theta}_t}(\tau) \right\} \right].$$
(A3)

Here,  $\tau=(a_1,a_2,\cdots,a_q)$  is the action sequence sampled from the previous policy  $\pi_{\theta_t}$  [cf. Algorithm 1]. Typically, the policy from the last iteration is chosen to be the old policy;  $\rho_{\theta}(\cdot) = \frac{\pi_{\theta}(\cdot)}{\pi_{\theta_t}(\cdot)}$  is the importance sampling weight between the new policy  $\pi_{\theta}$  and the old policy  $\pi_{\theta_t}$ ;  $A_{\theta_t}(\tau) = R(\tau) - b$  is the advantage function, where b is called a baseline: the advantage measures the reward gain of choosing a specific action, w.r.t. the baseline. For example, a simple baseline can be the average reward, e.g.,  $b = \mathbb{E}_{\tau \sim \pi_{\theta_t}}[R(\tau)]$ , and then the advantage measures how much better (or worse) an action is w.r.t. the average; in the numerical experiments, we use an exponential moving average [cf. App. A 5 for the details].

Further, the clip function.

$$clip(r, x, y) = \max (\min (r, x), y),$$

clips the value of r within the interval [x, y], which is used

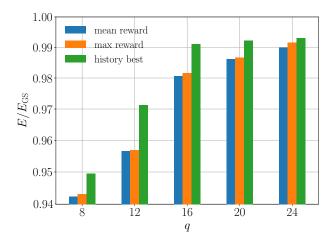


FIG. 12. Spin-1 Ising model: energy minimization against different circuit depths q using CD-QAOA. The mean negative energy density (blue) is computed for a sample generated using the final, learned policy; max (orange) is the maximum within the sample; the history best (green) is the best encountered policy during the entire training process (i.e., considering all iterations). The total duration T=20 and the values of q ranges from 8 to 24. The other model parameters are the same as in Fig. 11.

to restrict the likelihood ratio in the range  $[1-\epsilon, 1+\epsilon]$ ; this prevents the policy update from deviating too much from the old policy after one gradient update. The clipped objective function is designed to improve the policy as well as to keep it within some vicinity of the last iteration, whence the name Proximal Policy Optimization.

We update the network parameters  $\boldsymbol{\theta}$  by ascending along the gradient of the RL objective  $\mathcal{G}(\boldsymbol{\theta})$ . To provide intuition about the PPO objective, consider the following limiting case. If we only have the first term in the objective, i.e.  $\mathcal{G}_1(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}_t}}[\rho_{\boldsymbol{\theta}}(\tau)A_{\boldsymbol{\theta}_t}(\tau)]$ , we obtain the following gradient:

$$\nabla_{\boldsymbol{\theta}} \mathcal{G}_{1}(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}_{t}}} \left[ \nabla_{\boldsymbol{\theta}} \rho_{\boldsymbol{\theta}}(\tau) A_{\boldsymbol{\theta}_{t}}(\tau) \right]$$
$$= \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}_{t}}} \left[ \frac{\nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(\tau)}{\pi_{\boldsymbol{\theta}_{t}(\tau)}} A_{\boldsymbol{\theta}_{t}}(\tau) \right].$$

Since we are taking the gradient with respect to  $\boldsymbol{\theta}$ , it will pass through  $\pi_{\boldsymbol{\theta}_t}$  and  $A_{\boldsymbol{\theta}_t}(\tau)$ . Furthermore, whenever the parameters  $\boldsymbol{\theta} \approx \boldsymbol{\theta}_t$ , the gradient above is identical to the policy gradient:

$$\nabla_{\boldsymbol{\theta}} \mathcal{G}_{1}(\boldsymbol{\theta}) \approx \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}}} \left[ \frac{\nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(\tau)}{\pi_{\boldsymbol{\theta}(\tau)}} A_{\boldsymbol{\theta}}(\tau) \right]$$
$$= \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}}} [\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\tau) A_{\boldsymbol{\theta}}(\tau)].$$

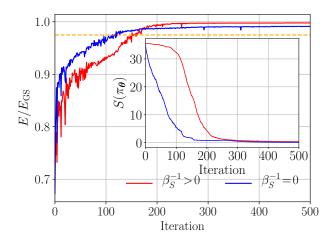


FIG. 13. Spin-1 Ising model: Comparison of the mean reward with  $(\beta_{S,\{0\}}^{-1}=0.1)$  and without  $(\beta_S^{-1}=0)$  the entropy bonus during training. For comparison, the dashed horizontal line marks the performance of QAOA. The inset shows the evolution of the policy information entropy during training. Adding entropy gives more room for the RL agent to explore the space of policies instead of directly exploiting the knowledge it obtains. As becomes clear from the figure, the RL algorithm with the entropy bonus achieves a better final performance at the end of training, at the cost of suffering an intermediate lower reward at the beginning of training. The simulation parameters are the same as in Fig. 11.

However, PPO performs multiple gradient updates on the sampled data, rendering policy learning more sample efficient [85].

# a. Incentivizing Exploration using Entropy

Maintaining a balance between exploration and exploitation is another major challenge for the reinforcement learning algorithm. Too much exploration prevents the agent from adopting the best strategy it knows so far; on the contrary, too much exploitation limits the agent from attempting new actions and achieving a potentially higher reward. Therefore, it is more appropriate for the agent to explore substantially in the initial iterations of the training procedure, and to gradually switch over to exploitation towards the end of the training procedure.

In order to incentivize the agent to explore the action space at the beginning of training, we include an entropy 'bonus' [86, 87] to the PPO objective from Eq. (A3). To do this, consider the maximal-entropy objective, where the agent aims to maximize the sum of the total reward and the policy entropy S [cf. Eq. (A5)]:

$$\mathcal{J}(\boldsymbol{\theta}) = \mathcal{G}(\boldsymbol{\theta}) + \beta_S^{-1} \mathcal{S}(\pi_{\boldsymbol{\theta}}) 
= \mathbb{E}_{\tau = (a_1, \dots, a_q) \sim \pi_{\boldsymbol{\theta}_t}} \left[ \min \{ \rho_{\boldsymbol{\theta}}(\tau) A_{\boldsymbol{\theta}_t}(\tau), \operatorname{clip} \left( \rho_{\boldsymbol{\theta}}(\tau), 1 - \epsilon, 1 + \epsilon \right) A_{\boldsymbol{\theta}_t}(\tau) \} + \beta_S^{-1} \sum_{j=1}^q \mathcal{S} \left( \pi_{\boldsymbol{\theta}}(\cdot | a_1, \dots, a_{j-1}) \right) \right], \tag{A4}$$

where  $\mathcal{S}(\pi_{\theta}(\cdot|a_1,\cdots,a_{j-1})) \equiv \mathcal{S}(\pi_{\theta}(\cdot))$ , for j=1. The trade-off between exploration and exploitation is controlled by the coefficient  $\beta_S^{-1}$ , which carries a meaning analogous to temperature in statistical mechanics: for  $\beta_S^{-1} \to 0$  (or  $\beta_S \to \infty$ ), any exploration is limited to the intrinsic probabilistic nature of the policy; if training is successful, it is expected that, for deterministic environments, the policy eventually converges to a delta distribution (over the action space) at the later training iterations; this may deteriorate exploration and learning. However, in the opposite limit,  $\beta_S^{-1} \to \infty$  (or  $\beta_S \to 0$ ), every action is selected with equal probability, and the values of the policy  $\pi$  become irrelevant. Therefore, in practice, we use a decay schedule for the inverse temperature  $\beta_S^{-1}$  to gradually reduce exploration [see App. A 5].

Since the marginal distribution  $\pi_{\theta}(\cdot)$  and the conditional distribution  $\pi_{\theta}(\cdot|a_1,\cdots,a_{j-1})$  are discrete categorical distributions over  $\mathcal{A}$ , we can compute a closed form expression for the entropy of the categorical distribution policy. For trajectory  $\tau^i = (a_1^i, \cdots, a_q^i)$ , the j-th term in the entropy bonus simplifies to

$$S(\pi_{\theta}(\cdot|a_1^i,\cdots,a_{j-1}^i))$$

$$= -\sum_{a \in \mathcal{A}} \pi_{\theta}(a|a_1^i,\cdots,a_{j-1}^i) \log \pi_{\theta}(a|a_1^i,\cdots,a_{j-1}^i).$$
(A5)

We emphasize that the entropy considered here is the Shannon or information entropy associated with the policy as a probability distribution, and should be contrasted with the thermodynamic entropy, associated with the logarithm of the density of protocol configurations (a.k.a. density of states) in the optimization landscape. The Shannon entropy helps exploration in the space of policies, and thus the annealing of the corresponding Lagrange multiplier,  $\beta_S^{-1}$ , is not related to thermal annealing in the optimization/energy landscape in a straightforward manner. Moreover, notice that the policy optimization is part of the classical postprocessing of the quantum data, i.e., it does not compromise the nature of the quantum data which is fed to the algorithm in form of rewards.

Figure 13 shows a comparison of PPO with and without entropy, as controlled by the value of the temperature  $\beta_S^{-1}$ . Introducing the policy information entropy keeps the policy a bit broader in the initial stages of training which enhances exploration; towards the end of training the information entropy is not needed: therefore, we gradually "anneal"  $\beta_S^{-1}$ , cf. App. A 5.

# 5. Technical Details

We train the CD-QAOA algorithm for 500 epochs/iterations with a mini-batch size of M=128. Throughout the training, we sample trajectories according to the marginal and conditional policy distributions given by the autoregressive network.

We use Adam to perform gradient descent on the objective in Eq. (A4) with the default parameters  $\beta_1=0.9$  and  $\beta_2=0.999$ , which define the exponential decay rate for the first and second moment estimates, respectively. The learning rate is initialized as  $\alpha_{\{lr,0\}}=0.01$  and decays by a factor of 0.96 every 50 steps in a staircase fashion. To be more precise, the learning rate at the k-th iteration with the exponential decay reads as  $\alpha_{lr,\{k\}}=0.01\cdot0.96^{\lfloor k/50\rfloor}$ . The subscript  $\{k\}$  denotes the iteration/episode number.

We also introduce an exponential decay schedule for the pre-factor [a.k.a. temperature]  $\beta_S^{-1}$  of the entropy bonus from Eq. (A4). The temperature initializes at  $\beta_{S,\{0\}}^{-1} = 0.1$  and decays by a factor of 0.9 every 10 steps. At the k-th iteration, the temperature is  $\beta_{S,\{k\}}^{-1} = 0.1 \cdot 0.9^{k/10}$ . Eventually, the temperature is annealed to zero

We estimate the advantage function by  $A_{\theta_{\mathrm{old}}}(\tau) = R(\tau) - b$ , where b is the baseline used to reduce the variance of the estimation. Our baseline b uses an exponential moving average (EMA) of the previous rewards. EMA stabilizes the training and also leverages the past reward information to form a lagged baseline. In practice, we find that the RL algorithm can achieve better rewards compared with using the average of current samples as the baseline. To be more specific, the exponential moving baseline update is  $b_{\{k\}} = \eta b_{\{k-1\}} + (1-\eta)\bar{R}_{\{k\}}$ , where  $b_{\{0\}} = 0$  and  $\eta = 0.95$ . Here,  $\bar{R}_{\{k\}}$  is the sample average of the reward at the k-th iteration, i.e.  $\bar{R}_{\{k\}} = \frac{1}{M} \sum_{i=1}^{M} R_{\{k\}}^{i}(\tau^{i})$ .

In terms of policy optimization, we perform multiple steps of ADAM on the objective [Eq. (A4)]. The gradient update steps are 4 per minibatch. The clipped parameter in the objective is set to  $\epsilon$ =0.1.

The hyperparameters of the algorithm are listed in Table III.

# **Algorithm 1** CD-QAOA with autoregressive network based policy

Input: batch size M, learning rate  $\eta_t$ , total number of iterations  $T_{\text{iter}}$ , exponential moving average coefficient m, entropy coefficient  $\beta_S^{-1}$ , PPO gradient steps K.

- 1: Generate and select the gauge potential sets A using Algo. 2.
- 2: Initialize the autoregressive network and initialize the moving average  $\hat{R}=0$ .
- for  $t = 1, ..., T_{\text{iter}}$  do
- Autoregrssively sample a batch of discrete actions of size M, denoted as B:

$$\tau^k = (a_1^k, a_2^k, \dots, a_q^k) \sim \pi_{\theta}(a_1, a_2, \dots, a_q), \ k = 1, 2, \dots, M.$$

Apply the SLSQP solver to the lower-level continuous optimization (cf. App B): 5:

$$\min_{\{\alpha_j^k\}_{j=1}^q} \ \left\{ N^{-1} E(\{\alpha_j^k\}_{j=1}^q, \tau^k) \bigg| \sum_{j=1}^q \alpha_j^k = T; \ 0 \leq \alpha_j^k \leq T \right\}.$$

6:

Use the negative energy density as the return and compute the moving average: 
$$R_k = -N^{-1}E(\{\alpha_j^k\}_{j=1}^q, \tau^k), \quad \hat{R} = m \cdot \hat{R} + (1-m) \cdot \frac{1}{M} \sum_{k=1}^M R_k.$$

- Compute the advantage estimates  $A_k = R_k \hat{R}$ . 7:
- Initialize the parameter  $\boldsymbol{\theta}_{t+1}^{[1]} = \boldsymbol{\theta}_t$ . 8:
- for  $\kappa = 1, ..., K$  do 9:
- Evaluate the likelihood of samples using the parameters from last iteration and current iteration, i.e.  $\pi_{\theta_t}(\tau^k)$ , 10:  $\pi_{\boldsymbol{\theta}_{t+1}^{[\kappa]}}(\tau^k)$ , and compute the importance weight  $\rho_k^{[\kappa]} = \pi_{\boldsymbol{\theta}_{t+1}^{[\kappa]}}(\tau^k)/\pi_{\boldsymbol{\theta}_t}(\tau^k)$ .
- Use the advantage estimate and importance weight to compute  $\mathcal{G}_k, \mathcal{S}_k$ , following Eq. (A3) and Eq. (A5). 11:
- Compute the CD-QAOA objective Eq. (A4) and backpropagate to get the gradients: 12:

$$\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}_{t+1}^{[\kappa]}) = \frac{1}{M} \sum_{\left\{a_{j}^{\left\{k\right\}}\right\}_{j=1}^{q} \in B} \nabla_{\boldsymbol{\theta}} \left[ \mathcal{G}_{k}^{[\kappa]} + \beta_{S}^{-1} \mathcal{S}_{k}^{[\kappa]} \right].$$

- Update weights  $\boldsymbol{\theta}_{t+1}^{[\kappa+1]} \leftarrow \boldsymbol{\theta}_{t+1}^{[\kappa]} + \eta_t \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}_{t+1}^{[\kappa]})$ . 13:
- Update the parameter  $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_{t+1}^{[K+1]}$ 14:

# Appendix B: Low-level optimization: finding optimal protocol time steps $\alpha_i$

In order to determine the values of the time steps  $\alpha_i$ , we proceed as follows. For any given sequence of actions (or protocol sequence)  $\tau = (a_1, \dots, a_q)$  of total duration T, we solve the following low-level optimization problem:

$$\min_{\{\alpha_j\}_{j=1}^q} \left\{ N^{-1} E(\{\alpha_j\}_{j=1}^q, \tau) \middle| \sum_{j=1}^q \alpha_j = T; \ 0 \le \alpha_j \le T \right\}$$
(B1)

where q is the sequence length (circuit depth), N is the system size, and  $E(\cdot)$  is the energy of the final quantum state [cf. Eq. (2)] after evolving the initial quantum state  $|\psi_i\rangle$  according to the fixed protocol  $\tau$ .

Note that the  $\alpha_i$ -optimization is both bounded and constrained. It fits naturally into the framework of the Sequential Least Squares Programming (SLSQP). SLSQP solves the nonlinear problem in Eq. (B1) iteratively, using the Han-Powell quasi-Newton method with a Broyden-Fletcher-Goldfarb-Shanno (BFGS)[89] update of the B-matrix (an approximation to the Hessian matrix), and an L1-test function within the step size.

During each iteration of the policy update, a batch of trajectories  $\{\tau^i\} = \{(a_1^i, \dots, a_q^i)\}_{i=1}^M$  is sampled. Each trajectory sequence  $\tau^i$  is assigned a reward, by solving the optimization problem in Eq. (B1). Since performing the low-level optimization in Eq. (B1) is independent of the high-level optimization discussed in App. A, we run the former concurrently to boost the efficiency of the algorithm. We distribute every sequence  $\tau^i = (a_1^i, \dots, a_q^i)$ to a different worker process and aggregate the results back to the master process in the end. In practice, we use the batch size M = 128, and we distribute the simulation on 4 nodes with 32 cores each, so that each core solves only one optimization at a time.

Recently, it was demonstrated that it is possible to perform the continuous optimization on par with the discrete one, which eliminates the need to use a solver and results in a fully RL optimization approach [90].

# Appendix C: Scaling with the number of particles N, the protocol duration T, and the circuit depth q

Next, we discuss the computational scaling of CD-QAOA. While there are a number of (hyper-)parameters

Parameter	Value
Farameter	varue
optimizer	Adam [88]
learning rate $(\eta_{\{0\}})$	$1 \cdot 10^{-2}$
learning rate decay steps	50
learning rate decay factor	0.96
learning rate decay style	Staircase
RL temperature $(\beta_{S,\{0\}}^{-1})$	$1 \cdot 10^{-1}$
RL temperature decay steps	10
RL temperature decay factor	0.9
RL temperature decay style	Smooth
baseline exponential moving decay factor $(m)$	0.95
gradient steps (PPO)	4
clip parameter $\epsilon$	0.1
number of hidden layers	2
number of hidden units per layer $(d_{\text{hidden}})$	112
nonlinearity	ReLU
number of samples per minibatch $(M)$	128

TABLE III. Hyperparameter values for training the autoregressive deep learning model. In the case of  $|\mathcal{A}_{\text{CD-QAOA}}| = 9$ , q = 18 [cf. Eq. (3)] the total number of parameters is 24431; for  $|\mathcal{A}_{\text{CD-QAOA}}| = 7$ , q = 20 [cf. Eq. (7)] the total number of parameters is 21985.

in the algorithm, here we focus on the system size N, the protocol duration T, and the circuit depth q – which are physically the most relevant ones. We also consider the continuous and discrete optimization steps separately (the continuous step being also an essential part of conventional QAOA).

When it comes to the continuous optimization performed by a solver [cf. App. B], the main computational cost comes from the quantum evolution itself. The basic operation inside the solver is a multiplication of the matrix exponential  $\exp(-i\alpha_j H_j)$  by the state  $|\psi_i\rangle$ . The Hamiltonian  $H_i$  is stored as a sparse matrix, and the action of the matrix exponential onto the quantum state,  $\exp(-i\alpha_i H_i) |\psi_i\rangle$ , can be evaluated without computing the matrix exponential itself with the help of a sparse matrix-vector product; this operation scales exponentially with the system size N, i.e.  $O(\exp(cN))$  for some constant c. If we denote the sequence length (a.k.a. circuit depth) by q, then the total cost for evaluating a single value of the continuous angle  $\alpha$  scales as  $O(q \exp(cN))$ . We stress that this cost is also incurred by conventional QAOA.

For the discrete optimization performed using reinforcement learning [App. A 4], notice first that the machine learning model is agnostic to the physical quantum model, because we do not use information about the quantum model to train the policy, cf. App. A 2. Because the policy input is, by construction, independent of the quantum state, the input layer of the neural network architecture is shielded from the exponential growth of the physical Hilbert space with N. Hence, the deep neural network is *independent* of the Hilbert space dimension. Further, we use an autoregressive network model which scales linearly with the sequence length q, and also linearly with the size of the available action set  $|\mathcal{A}|$ . Thus,

N	T	q	$t_{\rm solver} \; ({ m sec/iter})$	$t_{\rm RL}~({ m sec/iter})$
10	20	20	$57.254 \pm 13.829$	$0.042 \pm 0.005$
8	20	20	$17.24 \pm 2.554$	$0.055 \pm 0.024$
6	20	20	$10.559 \pm 3.963$	$0.028\pm0.004$
4	20	20	$6.021 \pm 5.149$	$0.027 \pm 0.002$
10	28	20	$68.55 \pm 19.044$	$0.055 \pm 0.019$
10	24	20	$61.425 \pm 15.171$	$0.038\pm0.009$
10	20	20	$57.254 \pm 13.829$	$0.042\pm0.005$
10	16	20	$49.043 \pm 12.447$	$0.041\pm0.007$
10	12	20	$39.33 \pm 13.976$	$0.038\pm0.006$
10	8	20	$24.689 \pm 14.348$	$0.033\pm0.008$
10	4	20	$7.023 \pm 2.651$	$0.025 \pm 0.001$
8	20	24	$20.723 \pm 3.903$	$0.065 \pm 0.024$
8	20	20	$17.24 \pm 2.554$	$0.055 \pm 0.024$
8	20	16	$12.626 \pm 3.129$	$0.024\pm0.004$
8	20	12	$8.641 \pm 2.654$	$0.02\pm0.003$
8	20	8	$5.511 \pm 2.18$	$0.016 \pm 0.002$
8	20	4	$2.092 \pm 1.312$	$0.011\pm0.002$

TABLE IV. Wall clock running time of the two-level CD-QAOA optimization steps for the with different system sizes N, protocol durations T, and circuit depths q. The right-hand side of the table shows the time used for the lower-level solver (column  $t_{\rm solver}$ ) and the time spent for the high-level RL algorithm (column  $t_{\rm RL}$ ) at every successful iteration. The total cost can then be obtained by multiplying the time for  $t_{\text{solver}}$  by the appropriate number of repetitions (e.g., continuous solver realizations, policy sample batch size, PPO training episodes, etc.), taking into account any parallelization if present. Every number represents an average over 40 independent runs with the corresponding standard deviation shown; the significant deviation in  $t_{\text{solver}}$  is caused by the random initial solver state used which causes the algorithm to take a different number of steps to converge within the given tolerance, cf. App. D. This test is carried out on a single processor Intel Core i7-8700K CPU 6-core 3.70GHz.

the total computational cost for the reinforcement learning optimization scales as  $O(q|\mathcal{A}|)$ . The scaling of the neural network with the variational networks parameters (weights and biases) is trivially given by the matrix-vector multiplication, as is the case for typical ML deep networks, and is also independent of the physics of the controlled system.

A comparison of the wall clock time for the discrete and continuous optimization steps is provided in Table IV. We distinguish between the continuous solver optimization and the discrete RL optimization, and show the average times for one successful step of each in the two columns on the right-hand-side. The total cost can then be obtained by multiplying the time for  $t_{\rm solver}$  by the appropriate number of repetitions (e.g., continuous solver initial conditions, policy sample batch size, PPO training episodes, etc.), and by multiplying the time for  $t_{\rm RL}$  by the

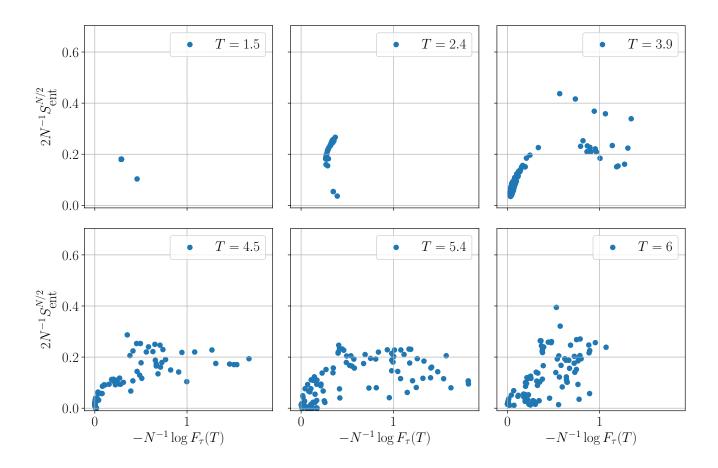


FIG. 14. Spin-1/2 Ising model: Visualization of the continuous optimization landscape for the durations  $\alpha_j$  in the fidelity-entanglement entropy plane, for the best sequence found by the RL agent [see App. D]. Each point corresponds to a local minimum, obtained using the SLSQP optimizer, starting from a uniformly drawn random initial condition. The system size is N = 16, and the rest of the parameters are the same as in Fig. 1.

number of PPO iterations, thereby taking into account any parallelization if used; for instance, the most expensive simulation we performed ran for about 109 hours on four nodes (Intel Xeon Skylake 6130 32-core 2.1 GHz) to produce the  $N=10,\,T=12,\,q=20$  data point shown in Table II.

We emphasize that the time  $t_{\rm solver}$  required for the continuous optimization is an essential part of conventional QAOA, and is the current limiting factor for reaching large system sizes, as is the case in merely all simulations of quantum dynamics on classical computing devices. In sharp contrast, the cost for training the deep autoregressive network is N-independent, and  $t_{\rm RL}$  per iteration is negligible; however, the choice of RL algorithm can strongly impact the number of iterations. CD-QAOA is, thus, suitably designed for potential applications on quantum simulators and quantum computers which will enable accessing large system sizes bypassing the exponential bottleneck intrinsic to simulations of quantum dynamics on classical devices.

# Appendix D: Many-Body Control Landscape

Let us briefly address the question about how hard the many-body ground state preparation problems are, that we introduced in the main text. To this end, recall that CD-QAOA has a two-level optimization structure: (i) discrete optimization to construct the optimal sequence of unitaries [App. A], and (ii) continuous optimization to find the best angles, given the sequence, to minimize the cost function [App. B]. Here, we focus exclusively on the continuous optimization landscape, and postpone the discrete landscape to a future study.

The RL agent learns in batches/samples of M=128 sequences, which sample the current policy at each iteration step and provide the data set for the policy gradient algorithm. To evaluate each sequence in the batch, we use SLSQP to optimize for the durations  $\alpha_j$  in a constrained and bounded fashion:  $\sum_j \alpha_j = T$  and  $0 \le \alpha_j \le T$  [cf. App. B]. This provides us with the full unitary  $U(\{\alpha_j\}_{j=1}^q, \tau)$ ; applying it to the initial state we obtain the reward value for this sequence. This procedure repeats iteratively as the RL agent progressively

discovers improved policies.

Once the RL agent has learned an optimal sequence, i.e. after the optimization procedure is complete, we focus on the best sequence from the sample, and examine how difficult it is to find the corresponding durations  $\alpha_j$  using SLSQP. To this end, we draw q values at random from a uniform distribution over the interval [0, T/q], and use them as initial conditions for the  $\alpha_j$ , to initialize the SLSQP optimizer with. We use the same q as the circuit depth so that the initial durations  $\alpha_j^{(0)}$  are, on average, equal. We then repeat this procedure P times, and generate a sample  $\mathcal{M} = \left\{ \left\{ \alpha_j^n \right\}_{j=1}^q \right\}_{n=1}^P$  of the local minima in the optimization landscape for  $\alpha_j$ 's. The larger P, the better our result for the true reward assigned to  $\tau$  is.

Notice that, in the beginning of the training, the RL agent is still in the exploration stage and the reward estimation does not need to be too accurate; this reward estimation needs to be more accurate as the agent switches over exploitation during the end of the training. In order to make the algorithm computationally more efficient, we introduce a linear schedule for the number of realizations of the  $\alpha_j$ -optimizer, starting from 3 with an increment of 1 every 30 iteration steps, i.e.  $P_{\{k\}}^{\rm tot} = 3 + \lfloor k/30 \rfloor$ , where subscript k indicates the iteration number for the RL policy optimization. In order to further save time in the reward estimation, we also introduce some randomness here by sampling  $P_{\{k\}}$  from a uniform distribution over  $1, 2, \dots, P_{\{k\}}^{\rm tot}$ .

Even though they all correspond to the same sequence, every local minimum in  $\mathcal{M}$  represents a potentially different protocol, since the durations  $\alpha_i$  will cause the initial quantum state to evolve into a different final state. We can evaluate for every protocol in  $\mathcal{M}$  the negative logfidelity,  $-\log F_{\tau}(T)$ , and entanglement entropy of the half chain,  $S_{\text{ent}}^{N/2}$ . Since the target state for the Ising model is an ordered ground state, it has area-law entanglement. Figure 14 shows a cut through the landscape in the fidelity-entanglement entropy plane for a few different durations T for the spin-1/2 Ising model. The better solutions are located in the lower left corner. The proliferation of local minima across the quantum speed limit has recently been studied in the context of RL [28] and QAOA [41]. This behavior indicates the importance of running many different SLSQP realizations, or else we may mis-evaluate the reward of a given sequence and the policy gradient will perform poorly.

Figure 14 also provides a plausible explanation for the destruction of the scaling collapse for  $T \gtrsim T_{\rm QSL}$  [Fig. 2]. Although the precision of the SLSQP optimizer is set at  $10^{-6}$ , the energy curves for large durations no longer fall on top of each other with a larger relative error. Hence, the occurrence of many local minima of roughly the same reward, which correspond to different protocols, effectively removes any universal features from the obtained solution; therefore, different system size simulations end up in different local minima.

### Appendix E: Variational Gauge Potentials

Consider the generic Hamiltonian

$$H(\lambda) = H_0 + \lambda H_1, \tag{E1}$$

with a general smooth function  $\lambda = \lambda(t)$ . We define a state preparation problem where the system is prepared in the ground state of  $H_0$  at time t = 0, and we want to transfer the state population in the ground state of H by time t = T.

Unlike adiabatic protocols, counter-diabatic driving relaxes the condition of being in the instantaneous ground state of  $H(\lambda)$  during the evolution. The idea is to reach the target state in a shorter duration T (compared to the adiabatic time) at the expense of creating controlled excitations [w.r.t. the instantaneous  $H(\lambda)$ ] during the evolution, which are removed before reaching the final time T. To achieve this, one can define a counter-diabatic Hamiltonian  $H_{\rm CD}$ . In general, the original  $H(\lambda)$  differs from  $H_{\rm CD}$ , whose ground state the system follows adiabatically:

$$H_{\rm CD}(\lambda) = H(\lambda) + \dot{\lambda} A_{\lambda},$$
 (E2)

where  $A_{\lambda}$  is the gauge potential;  $A_{\lambda}$  is defined implicitly as the solution to the equation [18]

$$[\partial_{\lambda} H + i[A_{\lambda}, H], H] = 0. \tag{E3}$$

The boundary conditions  $H_{\text{CD}}(\lambda(0)) = H(\lambda(0))$  and  $H_{\text{CD}}(\lambda(T)) = H(\lambda(T))$  impose the additional constraint  $\dot{\lambda}(0) = 0 = \dot{\lambda}(T)$  which suppresses excitations at the beginning and at the end of the protocol.

Using Eq. (E3), one can convince oneself that the gauge potential  $A_{\lambda}$  of a real-valued Hamiltonian H is always imaginary-valued [18].

For generic many-body systems, it has recently been argued that the gauge potential  $A_{\lambda}$  is a nonlocal operator [34]. Nevertheless, one can proceed by constructing a variational approximation  $\mathcal{X} \approx A_{\lambda}$ , which minimizes the action

$$S(\mathcal{X}) = \langle G^2(\mathcal{X}) \rangle - \langle G(\mathcal{X}) \rangle^2, \quad G(\mathcal{X}) = \partial_{\lambda} H + i[\mathcal{X}, H].$$
(E4)

For ground state preparation,  $\langle \cdot \rangle = \langle \psi_{\rm GS}(\lambda) | \cdot | \psi_{\rm GS}(\lambda) \rangle$  is the instantaneous ground state expectation value w.r.t.  $H(\lambda)$ . More generally, one can use  $\langle \cdot \rangle = \text{Tr}(\rho_{\rm th} \times (\cdot))$ , where  $\rho_{\rm th} \propto \exp(-\beta H)$  is a thermal density matrix at temperature  $\beta^{-1}$ : for  $\beta \to \infty$ , we recover the ground state expectation value; for  $\beta \to 0$  all eigenstates are weighted equally.

We mention in passing that alternative schemes to approximation the adiabatic gauge potential have also been considered [37].

#### 1. Spin Hamiltonians

### a. Real-valued Spin-1/2 Hamiltonians

Let H now be a real-valued spin-1/2 Hamiltonian with translation and reflection invariance. Such a system is given, e.g., by the mixed-field Ising model, discussed in the main text. We now construct an ansatz for the variational gauge potential  $\mathcal X$  which obeys these symmetries, and is imaginary valued.

We can organize the terms contained in  $\mathcal{X}$  according to their multi-body interaction type, as follows. The only single-body imaginary valued term we can write is  $\sum_j \beta_j S_j^y$ . Translation and reflection symmetries, whenever present in H, further impose that the coupling constant  $\beta_j = \beta$  be site-independent, i.e. spatially uniform. Hence, this is the zeroth-order term in our variational gauge potential construction, cf. Eq. (E5).

Next, we focus on the two-body terms. Because the exact  $A_{\lambda}$  is imaginary valued for real-valued Hamiltonians, we may only consider interaction terms where  $S^{y}$  appears precisely once:  $S^{x}S^{y}$  and  $S^{y}S^{z}$ . For spin-1/2 systems, the two operators have to act on different sites, or else one can further simplify their product to single-body operators using the algebra for Pauli matrices. Once again, translation invariance dictates that the coupling constants are uniform in space, while reflection invariance requires us to take a symmetric combination. Imposing further that the interaction be short-range (we want to

construct the most local variational ansatz), we arrive at

$$\mathcal{X}(\{\beta_{l}^{(k)}\}) = \sum_{j} \beta_{0}^{(0)}(\lambda) S_{j}^{y} + \beta_{1}^{(0)}(\lambda) \left(S_{j+1}^{x} S_{j}^{y} + S_{j+1}^{y} S_{j}^{x}\right) + \beta_{1}^{(1)}(\lambda) \left(S_{j+1}^{z} S_{j}^{y} + S_{j+1}^{y} S_{j}^{z}\right). \tag{E5}$$

The coefficients  $\beta_l^{(k)}$  are the variational parameters that we need to determine to find the approximate CD protocol. To find their optimal values, we minimize the action  $\mathcal{S}(\mathcal{X})$  [18]. Note that, since we do not have a closed-form expression for the instantaneous ground state of  $H(\lambda)$ , we do the minimization numerically at every fixed time t along the protocol  $\lambda(t)$  [cf. App. E 2].

We can, in principle, add the next order terms to the series; however, they will either be less local, or consist of three- and higher-body interactions, which is hard to implement in experiments.

#### b. Real-valued Spin-1 Hamiltonians

The situation is more interesting for spin-1 systems: the eight-dimensional Lie algebra  $\mathfrak{su}(3)$ , which generates SU(3), contains three distinct imaginary-valued directions, which form a closed subalgebra  $\mathfrak{su}(2) \subsetneq \mathfrak{su}(3)$ , and hence there is more room to generate imaginary-valued combinations. To find all imaginary-valued terms consistent with a set of symmetries, we use QuSpin's functionality to implement an algorithm [App. E3] that lists them for generic bases [54, 55].

Translation and Reflection Symmetric spin-1 Hamiltonians, such as the spin-1 Ising and Heisenberg models, have a similar expansion to their spin-1/2 counterparts, but allow for more terms. Restricting the expansion to two-body terms, we have

$$\mathcal{X}(\{\beta_{l}^{(k)}\}) = \sum_{j} \left[ \beta_{0}^{(0)}(\lambda) S_{j}^{y} + \beta_{1}^{(0)}(\lambda) \left( S_{j}^{x} S_{j}^{y} + S_{j}^{y} S_{j}^{x} \right) + \beta_{2}^{(0)}(\lambda) \left( S_{j}^{z} S_{j}^{y} + S_{j}^{y} S_{j}^{z} \right) + \right. \\ \left. + \beta_{0}^{(1)}(\lambda) \left( \left[ S_{j+1}^{x} - a S_{j}^{x} \right] S_{j}^{y} + \left[ S_{j+1}^{y} - a S_{j}^{y} \right] S_{j}^{x} \right) + \beta_{1}^{(1)}(\lambda) \left( \left[ S_{j+1}^{z} - b S_{j}^{z} \right] S_{j}^{y} + \left[ S_{j+1}^{y} - b S_{j}^{y} \right] S_{j}^{z} \right) \right].$$
(E6)

where the constants a and b are chosen so that all five terms are mutually orthogonal w.r.t. the scalar product induced by the trace (i.e. Hilbert-Schmidt) norm; this ensures the linear independence of the constituent terms. Note that the three imaginary-valued on-site terms correspond precisely to the imaginary-valued  $\mathfrak{su}(2) \subseteq \mathfrak{su}(3)$ .

Adding Magnetization Conservation and Spin Inversion Symmetry further reduces the allowed terms in the series. Therefore, one has to restrict to three- and four-body terms:

$$\mathcal{X}(\{\zeta_{l}^{(k)}\}) = \sum_{j} \zeta_{0}^{(2)}(\lambda) \left(iS_{j}^{+}S_{j+1}^{-}S_{j+2}^{z} + iS_{j}^{z}S_{j+1}^{-}S_{j+2}^{+} + \text{h.c.}\right) + \zeta_{0}^{(3)}(\lambda) \left(iS_{j}^{-}S_{j}^{z}S_{j+1}^{+}S_{j+1}^{z} + iS_{j}^{+}S_{j}^{z}S_{j+1}^{-}S_{j+1}^{z} + \text{h.c.}\right),$$
(E7)

Because these terms are multi-body and less local, we re-

frain from using them in CD-QAOA in the present study.

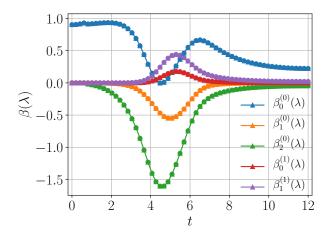


FIG. 15. Spin-1 Ising chain. Time dependence of the optimal coefficients  $\beta_l^{(k)}(\lambda(t))$  in the variational gauge potential (Eq. E6) with translation and reflection symmetry, determined from the procedure in App. E 2. The total duration T=12 with the time discretization step  $\Delta t=0.2$ , and the system size N=8. The protocol we used is  $\lambda(t)=\sin^2\left(\frac{\pi t}{2T}\right)$ . The other model parameters are the same as in Fig. 7.

We merely list them here for completeness.

As explained in the main text, to apply CD-QAOA for many-body ground state preparation, we consider the constituent terms in  $\mathcal{X}$  as independent generators  $\{H_j\}_{j=1}^{|\mathcal{A}|}$ . This comes in contrast to the variational gauge potential method where the ratios between the coefficients  $\beta_l^{(k)}$  play an important role.

# c. Variational Gauge Potential Ansatz for the Lipkin-Meshkov-Glick Model

As explained in the main text, the Lipkin-Meshkov-Glick (LMG) Hamiltonian, cf. Eq. (5), models homogeneously interacting spin-1/2 particles on an all-to-all connected graph in the presence of an external field. Here,

we compute the lowest-order terms appearing in the series for the variational gauge potential  $\mathcal{X}$ , going beyond Ref. [91].

The starting point is the LMG Hamiltonian

$$H = -\frac{J}{N} (S^x)^2 + h (S^z + N/2).$$
 (E8)

We introduce two bosonic modes, s and t, where  $S^z = t^{\dagger}t - N/2 = n_t - N/2$  and  $S^+ = t^{\dagger}s$ , and cast the LMG Hamiltonian in the form

$$H = ht^{\dagger}t - \frac{J}{4N} \left( t^{\dagger}s + s^{\dagger}t \right)^{2}. \tag{E9}$$

Recalling once again that real-valued Hamiltonians have imaginary-valued gauge potentials, and that gauge potentials do not have diagonal matrix elements, we make the following ansatz:

$$\mathcal{X}(\{\beta_l^{(k)}\}) = \beta_0^{(0)}(\lambda)Y + \beta_1^{(1)}(\lambda)\hat{XY} + \beta_1^{(0)}(\lambda)\hat{ZY}, \text{ (E10)}$$

where

$$Y = S^{y} = \frac{i}{2} \left( s^{\dagger} t - t^{\dagger} s \right),$$

$$\hat{X}Y = \frac{1}{N} \left( S^{x} S^{y} + S^{y} S^{x} \right) = -\frac{i}{2N} \left[ (t^{\dagger} s)^{2} - (s^{\dagger} t)^{2} \right],$$

$$\hat{Z}Y = \frac{1}{N} \left( \left( S^{z} + \frac{N}{2} \right) S^{y} + S^{y} \left( S^{z} + \frac{N}{2} \right) \right)$$

$$= \frac{i}{2N} \left( s^{\dagger} t^{\dagger} t t - s t^{\dagger} t t^{\dagger} t - s t^{\dagger} t^{\dagger} t - s t^{\dagger} t^{\dagger} t \right). \tag{E11}$$

To compute the matrix elements of the gauge potentials, we define the basis

$$|N, n_t\rangle = \frac{\left(t^{\dagger}\right)^{n_t} \left(s^{\dagger}\right)^{N-n_t}}{\sqrt{n_t!(N-n_t)!}}|0\rangle, \text{ with } n_t = 0, \dots, N.$$

The gauge potentials have the following non-zero matrix elements (plus their conjugates to make the operators hermitian):

$$\langle N, n_t | Y | N, n_t + 1 \rangle = -\frac{i}{2} \sqrt{(n_t + 1)(N - n_t)},$$

$$\langle N, n_t | \hat{XY} | N, n_t + 2 \rangle = \frac{i}{2N} \sqrt{(n_t + 2)(n_t + 1)(N - n_t - 1)(N - n_t)},$$

$$\langle N, n_t | \hat{ZY} | N, n_t + 1 \rangle = \frac{i}{2N} (2n_t + 1) \sqrt{(n_t + 1)(N - n_t)}.$$
(E12)

#### 2. Numerical Minimization to obtain the Variational CD Protocol

Since the action S in Eq. (E4) is quadratic in the variational parameters  $\beta_j$ , it is possible to derive a generic lin-

ear system, whose solutions are the optimal parameters of the variational gauge potential within CD driving [92].

Suppose that  $\mathcal{X} = \sum_{j=1}^{r} \beta_j H_j$  is given by a linear combination of r gauge potential terms. Then, it is straight-

forward to see that

$$G(\mathcal{X}) = \partial_{\lambda} H + \sum_{j=1}^{r} i[H_j, H] \beta_j.$$
 (E13)

Defining the operator-valued quantities  $B_0 = \partial_{\lambda} H$  and  $B_j = i[H_j, H]$  and setting  $\beta_0 = 1$ , we arrive at the following expression for the variational action

$$S(\mathcal{X}) = \left\langle \left( B_0 + \sum_j B_j \beta_j \right)^2 \right\rangle - \left( \langle B_0 + \sum_j B_j \beta_j \rangle \right)^2$$
$$= \sum_{i,j=0}^r \left( \langle B_i B_j \rangle - \langle B_i \rangle \langle B_j \rangle \right) \beta_i \beta_j, \tag{E14}$$

which is a quadratic form in the unknown coefficients  $\beta_j$ . To find the minimum of  $\mathcal{S}(\mathcal{X})$  w.r.t.  $\beta_j$ , we can take the derivative and set it to zero, to obtain the linear system of equations for the optimal  $\beta_j$ :

$$\sum_{k} \mathcal{M}_{jk} \beta_k = -\mathcal{M}_{0j} \tag{E15}$$

where  $\mathcal{M}_{jk} = \langle B_j B_k \rangle + \langle B_k B_j \rangle - 2 \langle B_j \rangle \langle B_k \rangle$ . Solving the system we obtain the minimum  $\{\beta_j\}_{j=1}^r$  of the variational action  $\mathcal{S}$ .

The ground state expectation values in the above procedure, as well as the Hamiltonian  $H(\lambda(t))$  depend implicitly on time  $t \in [0,T]$  via the protocol  $\lambda(t)$ . Therefore, to find the time dependence of  $\beta_j(t)$ , we discretize the time interval [0,T] into  $N_T$  time steps, and repeat the procedure at every time step. This yields  $\beta_j(t_i)$  at the time steps  $t_i$ . To recover the full functional dependence, we use a fine discretization mesh, and apply a linear interpolation to  $\beta_j(t_i)$ . Alternatively, notice that the coefficients  $\beta_j = \beta_j(\lambda(t))$  depend on time t only implicitly via the protocol  $\lambda$ . Therefore, it is also possible to discretize the range of  $\lambda(t)$  instead.

For the spin-1 Ising model, the time-dependence of  $\beta_j$  is shown in Fig. 15. This defines  $H_{\rm CD}$  which generates the CD evolution. In Sec. V and App. F 4, we compare variational CD driving to CD-QAOA and conventional QAOA.

# 3. Algorithm for Generating Gauge Potential Terms in the Presence of Lattice Symmetries

Finally, we also show the algorithm we used to determine the terms appearing in the gauge potential expansions in Eqn. (E5), (E6), and (E7), which obey a fixed set of symmetries.

In general, one can represent any local operator of the kind  $J_{i_1,\dots,i_l}O_{i_1}^{\gamma_1}\cdots O_{i_l}^{\gamma_l}$  as a triple  $(\mathcal{Y}, \mathcal{I}, J)$ , where  $J = J_{i_1,\dots,i_l}$  is the coupling coefficient constant,  $\mathcal{I} = (i_1,\dots,i_l)$  is the set of sites the operators act on, and  $\mathcal{Y} = (\gamma_1,\dots,\gamma_l)$  defines the types of operators that act

on the corresponding sites; the triple  $(\mathcal{Y}, \mathcal{I}, J)$  can then be used to construct the operator.

In the following, we refer to the separate terms appearing in the gauge potential series as 'Hamiltonians'  $H_j$ , i.e.  $\mathcal{X} = \sum_j \beta_j H_j$ ; a Hamiltonian is defined as  $H = \sum_{(i_1, \cdots, i_l) \in \Lambda} J_{i_1, \cdots, i_l} O_{i_1}^{\gamma_1} \cdots O_{i_l}^{\gamma_l}$ , where  $\Lambda$  is the lattice graph. As we argued above, real-valued Hamiltonians have purely imaginary-valued gauge potentials; thus, the coefficient J is chosen to be purely imaginary.

We build the series for the variational gauge potential  $\mathcal{X}$  recursively: we first consider a set  $\mathcal{L}_{\text{elem}}$  of elementary operators O — the building blocks for the expansion: e.g., for the spin-1 chains, these can be the spin-1 operators  $\mathcal{L}_{\text{elem}} = \{S^+, S^-, S^z\}$ . We want to construct the terms in the expansion for  $\mathcal{X}$  iteratively at a fixed order l, e.g. l=1 comprises single-body terms, l=2 — two-body terms, etc. We also assume that we have access to a routine which checks if a trial list of operators obeys a given lattice symmetry; if not, the same routine returns the missing operators to be added to the original list, so that the symmetry is now satisfied [e.g., such a routine is used in QuSpin [54, 55]].

The pseudocode we developed is shown in Algorithm 2. To construct multi-body terms at a fixed order l, we define combinations of the elementary operators, and store them in the list  $\mathcal{L}_{op}$ ; the way these combinations are built can be used to implement constraints, such as particle/magnetization conservation, etc. This is implemented via the product operator (Line 2 of Algorithm 2). It generates all possible combinations of selecting l elementary operators with replacement. The sets of lattice sites that the operators from  $\mathcal{L}_{op}$  act on, are stored in the list  $\mathcal{L}_{\text{sites}}$  (Line 3 of Algorithm 2). Then, for each trial triple  $(\mathcal{Y}, \mathcal{I}, J)$ , we make use of the routine to check the symmetry and record any operators which do not respect it. We append these, so-called *missing operators*, to the original list, and we keep checking the symmetry condition until we obtain all operators that satisfy the symmetry (Line 10-15 of Algorithm 2). The finite number of combinations guarantees a termination in a finite number of steps.

# Algorithm 2 Generation of variational gauge potential

**Input:** a list of required symmetries  $\mathcal{L}_{\text{sym}}$ , order l, a list of elementary operator types  $\mathcal{L}_{\text{elem}}$ .

- 1: Initialize empty list for gauge potential terms  $\mathcal{L}_{\text{gauge}}$ .
- 2: Generate all possible combinations of local operators at order  $\boldsymbol{l}$

$$\mathcal{L}_{\text{op}} = \text{product}(\mathcal{L}_{\text{elem}}, \text{repeat} = l).$$

3: Enumerate all possible combinations of lattice sites  $\mathcal{L}_{\text{sites}}$  the l-th order operators act on.

```
4: for \mathcal{Y} in \mathcal{L}_{op} do
          for \mathcal{I} in \mathcal{L}_{\mathrm{sites}} do
 5:
               Initialize an empty list \mathcal{L}_H
 6:
 7:
               Set J = i \ (i = \sqrt{-1}).
               Append (\mathcal{Y}, \mathcal{I}, J) to \mathcal{L}_H.
 8:
               Set the flag IsSym = False.
 9:
10:
               while IsSym is False do
                    Set IsSym = True.
11:
12:
                    for sym in \mathcal{L}_{sym} do
                         if exists missing operator (\mathcal{Y}', \mathcal{I}', J') then
13:
                              Set IsSym = False.
14:
                              Append (\mathcal{Y}', \mathcal{I}', J') to \mathcal{L}_H.
15:
               Build Hamiltonian H using the triplets in \mathcal{L}_H.
16:
               if H or equivalents not included in \mathcal{L}_{gauge} then
17:
                    Append H to \mathcal{L}_{\text{gauge}} .
18:
19: Return the list of gauge potential basis \mathcal{L}_{gauge}.
```

product: Cartesian product, equivalents: equivalent mod scalar, missing operator: the operator missed for the symmetry requirement

In order to avoid repeating previously identified Hamiltonians, we discard equivalent Hamiltonians (Line 17 of Algorithm 2): two Hamiltonians are called equivalent when one is a scalar times the other. Since here we consider imaginary-valued gauge potentials, the multiple constant should be real. To test whether the Hamiltonians  $H_1$  and  $H_2$  are equivalent in practice, it suffices to test whether  $H_1$  is equal to  $\pm \frac{\|H_1\|}{\|H_2\|} H_2$ , where we use the Hilbert-Schmidt norm.

# Appendix F: CD-QAOA for Many-Body State Preparation

Here, we provide a supplementary discussion on the performance of CD-QAOA for many-body pure state preparation using the quantum spin chains introduced in the main text. We refer the reader to the main text for the definition of various model parameters; the shorthand spin operator notation used is defined in Table I.

#### 1. Spin-1/2 Ising Chain

First, we show results for the single-spin problem (J = 0):

$$H = H_1 + H_2, H_1 = h_z S^z, H_2 = h_x S^x.$$
 (F1)

In Fig. 16, we clearly see that CD-QAOA [red curve] has a smaller quantum speed limit  $T_{\rm QSL} \approx 4.0$  than conventional QAOA [blue]; this is anticipated, since CD-QAOA

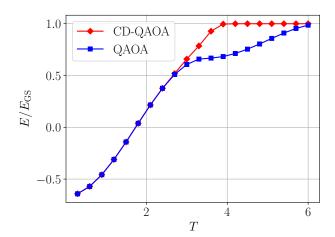


FIG. 16. Single spin-1/2 state preparation: energy density against protocol duration for CD-QAOA with  $\mathcal{A}_{\text{CD-QAOA}} = \{Z, X, Y\}$  (red) and conventional QAOA with  $\mathcal{A}_{\text{QAOA}} = \{Z, X\}$  (blue). The values of q is 3 for both methods. For conventional QAOA, we trained two possible alternating patterns (i.e.  $(Z \to X \to Z)$  and  $(X \to Z \to X)$ ) and pick the best one for the comparison. The model parameters are the same as in Fig. 1 with J=0.

has a larger control space at its disposal. Moreover, we find that, for  $T < T_{\rm QSL}$ , CD-QAOA only makes use of a single Y rotation by setting the durations  $\alpha_j$  associated with any other unitaries from the set  $\mathcal{A}$ , to zero. As mentioned in the main text, conventional QAOA tries to represent this Y-rotation by means of Euler angles, i.e. composed of X and Z rotations; in general, this results in a higher duration cost to complete the population transfer (leading to a larger  $T_{\rm QSL}$ ). However, for short durations T, a Y-rotation can be exactly obtained using a proper sequence of the X and Z terms. For these reasons, we find an exact agreement between the two curves for small values of  $T\lesssim 3$ .

Let us now switch on the spin-spin interaction strength J > 0; consider the spin-1/2 Ising chain

$$H = H_1 + H_2,$$

$$H_1 = \sum_{j=1}^{N} J S_{j+1}^z S_j^z + h_z S_j^z, \quad H_2 = \sum_{j=1}^{N} h_x S_j^x.$$
(F2)

Figure 17 [top] shows a comparison of the best learned energies, between conventional QAOA, and CD-QAOA for two sets  $(\mathcal{A}, \mathcal{A}')$  with different number of unitaries:  $|\mathcal{A}| = 5, |\mathcal{A}'| = 3$  [see caption]. We find that additionally using only the single-particle gauge potential term Y [green line], typically accessible in experiments, one can already obtain a higher-fidelity protocol than QAOA to prepare the ground state. Interestingly, for short protocol durations T, the two-body gauge potential terms, present in  $\mathcal{A}$  but not in  $\mathcal{A}'$ , do not contribute to improving the energy of the final state, as can be seen from the agreement of the red and green lines for  $T \lesssim 1.5$ . This suggests

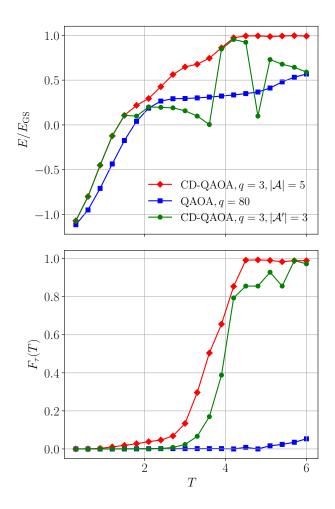


FIG. 17. Spin-1/2 Ising model: energy minimization (top) and many-body fidelity maximization (bottom) against protocol duration T. We compare CD-QAOA with  $\mathcal{A}_{\text{CD-QAOA}} = \{Z|Z+Z,X;Y,X|Y,Y|Z\}$  (red), CD-QAOA with  $\mathcal{A}'_{\text{CD-QAOA}} = \{Z|Z+Z,X;Y\}$  (green), and conventional QAOA with  $\mathcal{A}_{\text{QAOA}} = \{Z|Z+Z,X\}$  (blue). The model parameters are the same as in Fig. 1 with the number of spins N=14.

that single-particle processes dominate over many-body processes when it comes to lowering the energy of the z-polarized initial state, and implies that the target ground state is single-particle-like (i.e. close to a product state). The non-smooth behavior of the green curve at larger durations, is attributed to the ruggedness of the control landscape, as different runs of the SLSQP optimizer may get stuck in one of the many suboptimal local minima [App. D].

One may wonder if it is possible to prepare the ground state by straightforward fidelity maximization. We define the many-body fidelity to transfer the population to the target state using the unitary process  $U(\{\alpha_j\}_{j=1}^q, \tau)$ ,

with 
$$\sum_{j=1}^{q} \alpha_j = T$$
, as

$$F_{\tau}(T) = F(\{\alpha_j\}_{j=1}^q, \tau) = |\langle \psi_* | U(\{\alpha_j\}_{j=1}^q, \tau) | \psi_i \rangle|^2.$$
(F3)

The fidelity can be less relevant from the perspective of many-body physics because (i) the many-body fidelity is typically exponentially suppressed, and (ii) it requires a reference to the ground state itself (which we seek) in order to be computed. However, the fidelity of a quantum process is a widely used benchmark in quantum computing; it also provides a better measure (than energy density) for the distance between two states in the Hilbert space  $\mathcal{H}$ .

Figure 17 [bottom] shows the many-body fidelity for N=14 spins. Unlike the inset of Fig. 1 from the main text (where we show the fidelity associated with the protocol obtained using energy density minimization), here we use the fidelity as a reward function for QAOA. We observe that optimizing the fidelity behaves quantitatively similar to optimizing the energy density. We would like to emphasize here once again the advantage of the gauge potential ansatz: the conventional QAOA simulation is done using q=80 variational parameters  $\alpha_j$  [yet no significant improvement is observed for  $q\geq 4$ , cf. Fig. 1], while CD-QAOA requires only q=3 variational parameters.

Although the fidelity  $F_{\tau}(T)$  is anticipated to vanish for  $T < T_{\rm QSL}$  in the thermodynamic limit, the negative log-fidelity density,  $-N^{-1}\log F_{\tau}(T)$ , is more likely to. Figure 18 [inset] shows the finite size scaling of the fidelity curves. Similar to the energy density [Fig. 2], we obtain an almost perfect scale collapse. We verified that maximizing the fidelity produces similar results as minimizing the negative log-fidelity density for the spin-1/2 chain: at first sight, this is nontrivial because  $F_{\tau}(T)$  is exponentially suppressed with the system size N for  $T < T_{\rm QSL}$ ; however, this behavior is likely explained by the generalization capabilities of the RL agent from small to large system sizes [cf. Sec. VI].

# 2. Anisotropic Spin-1 Heisenberg Chain

Next, we discuss in detail the ground state preparation process in the anisotropic Heisenberg spin-1 chain:

$$H = H_1 + H_2,$$

$$H_1 = J \sum_{j=1}^{N} (S_{j+1}^x S_j^x + S_{j+1}^y S_j^y), \quad H_2 = \Delta \sum_{j=1}^{N} S_{j+1}^z S_j^z,$$
(F4)

where the model parameters are defined in the main text. An important detail worth mentioning is that the ferromagnetic ground state at  $\Delta/J=-2.0$  is two-fold degenerated (one state, corresponding to one of the two z-polarizations). While being a trivial observation, this requires certain care when analyzing the physics of the protocols the agent found. In particular, notice that energy minimization is insensitive to this degeneracy, and hence

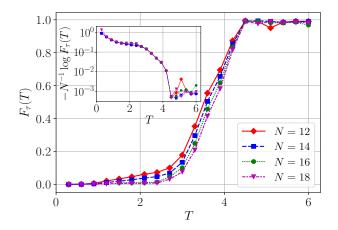


FIG. 18. Spin-1/2 Ising model: many-body fidelity maximization and corresponding quantity [inset, log scale] against protocol duration T for different system sizes N. The QAOA parameters are q=3 and  $\mathcal{A}=\{Z|Z+Z,X;Y,X|Y,Y|Z\}$ . The model parameters are the same as in Fig. 1.

the final state can appear as an arbitrary superposition of the two ferromagnetic states, and still have the correct ground-state energy. This leads to ambiguity when computing the fidelity of being in the target state: related to this, the cost function landscape likely develops a continuous one-dimensional structure for the (degenerate) global minima. Because we are interested in energy minimization, here we define the fidelity using the projector to the ground state manifold  $P = |\psi_*^{(1)}\rangle\langle\psi_*^{(1)}| + |\psi_*^{(2)}\rangle\langle\psi_*^{(2)}|$ :

$$F_{\tau}(T) = F(\{\alpha_j\}_{j=1}^q, \tau) = |\langle \psi_*^{(1)} | U(\{\alpha_j\}_{j=1}^q, \tau) | \psi_i \rangle|^2 + |\langle \psi_*^{(2)} | U(\{\alpha_j\}_{j=1}^q, \tau) | \psi_i \rangle|^2$$

where  $|\psi_*^{(1)}\rangle, |\psi_*^{(2)}\rangle$  are any two orthonormal states which span the doubly degenerate ground state manifold (e.g., the two FM ground states).

Figure 19 shows a comparison between CD-QAOA and conventional QAOA for FM, XY, and Haldane target states: the top row shows the result of energy density minimization [cf. Fig. 3]. The bottom row, on the other hand, displays the many-body fidelity associated with the same protocols. For  $\Delta/J=0.5$ , CD-QAOA allows reaching the target topological Haldane state already faster, as compared to conventional QAOA. Notice also that the gauge potential ansatz appears essential for reaching the target for both the XY ( $\Delta/J = -0.5$ ) and FM states  $(\Delta/J = -2.0)$ ; this becomes particularly obvious from the many-body fidelity curves. The latter also reveals an interesting detail: at  $\Delta/J = 0.5$ , a regime emerges around  $T \approx 5$ , where the QAOA fidelity is better than the CD-QAOA fidelity. However, this peculiarity below the quantum speed limit can be explained, recalling that the RL agent is given the (negative) energy density as the reward signal, and not the fidelity (note that CD-QAOA does outperform QAOA in energy).

In order to investigate in detail in the protocols found

by CD-QAOA, we fix a duration T, and consider the time evolution of the state,  $|\psi(t)\rangle = U(\{\alpha_j\}_{j=1}^q, \tau)|\psi_i\rangle$ , for three physical quantities:

(i) the energy

$$E(t) = \langle \psi(t) | H_* | \psi(t) \rangle$$

provides a measure of how far away in the cost function landscape the state is, at any given time  $t \in [0, T]$ .

(ii) the instantaneous fidelity

$$F_{\tau}(t) = |\langle \psi_* | \psi(t) \rangle|^2$$

(and its generalization to the doubly-degenerate ground state manifold), measures how far the current state is, from the target state  $|\psi_*\rangle$  in the Hilbert space; typically, we choose the ground state as the target state  $|\psi_*\rangle = |\psi_{\rm GS}(H)\rangle$ .

(iii) the entanglement entropy of the half chain

$$S_{\mathrm{ent}}^{N/2}(t) = -\mathrm{tr}_A \left[ \rho_A(t) \log \rho_A(t) \right], \ \rho_A(t) = \mathrm{tr}_{\bar{A}} |\psi(t)\rangle \langle \psi(t)|,$$

where A denotes a contiguous spacial region with a complement  $\bar{A}$  comprising half the periodic chain, and  $\rho_A(t)$  is the reduced density matrix on A at time t. For many-body systems, it is common to look at the entanglement entropy per site, which for spin-1 systems lies within the interval  $2N^{-1}S_{\rm ent}^{N/2} \in [0, \log 3]$ .

Figure 20 shows the time evolution of the energy, fidelity and entropy density, for all three target states of interest. For  $\Delta/J = 0.5$ , transferring the population from the AFM initial state to the Haldane state can be obtained equally well using either QAOA or CD-QAOA. Table V(d) shows the optimal protocol found by the RL agent: notice the three vanishing durations  $\alpha_2 = \alpha_{17} = \alpha_{18} = 0$ ; factoring them out, we recover precisely the conventional QAOA sequence (albeit with q odd). Thus, we see that the CD-QAOA may converge to conventional QAOA whenever the latter provides a highreward sequence. This result exemplifies our claim that CD-QAOA generalizes QAOA successfully. Of course, it is not clear whether this is the true global minimum of the cost function landscape (the RL agent does make use of the additional gauge potential terms for T < 7). Nevertheless, all physical quantities are expected to be prepared with similar accuracy under both protocols: to see this, notice that the entanglement entropy density depends only on the quantum state (unlike expectation values of observables), and that its value at t = T is close to the value for the target state (dashed horizontal line). Importantly, the entanglement remains area-law (as seen by the values being much smaller than the maximum entropy per site,  $\log(3)$ , suggesting the existence of a local effective Hamiltonian which generates the population transfer process dynamically.

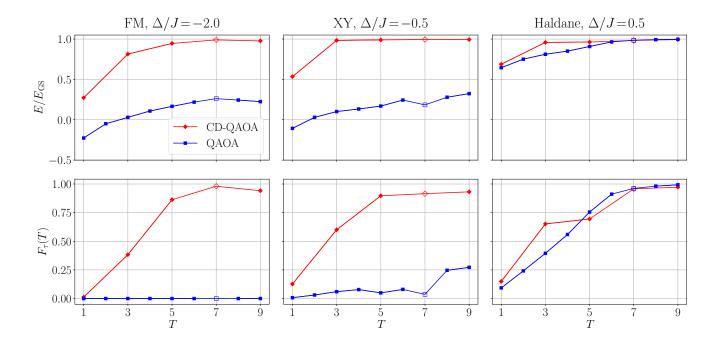


FIG. 19. Anisotropic Heisenberg spin-1 chain: energy minimization against protocol duration T — the corresponding energy (top row) and many-body fidelity (bottom row) for three ordered target states, corresponding to the ground state of the ferromagnetic (left,  $\Delta/J = -2.0$ ), XY (middle,  $\Delta/J = -0.5$ ), and Haldane (right,  $\Delta/J = 0.5$ ) target states, respectively. The empty symbols mark the duration at which we show the evolution of the system in Fig. 20. The model parameters are the same as in Fig. 3.

The best sequence for targeting the XY state at  $\Delta/J =$ -0.5 is shown in Table V(c). Although its structure is more complicated, factoring out the vanishing  $\alpha_i$ , we can discern two clear patterns: (i) the sequence starts and ends with two different single-particle basis rotations, and (ii) there is an alternating subsequence based on the subset  $\{X|X+Y|Y,Y\} \subseteq \mathcal{A}_{CD-QAOA}$ . Interestingly, the only gauge potential term used by the RL agent is the experimentally accessible single-particle Y rotation, and it is sufficient to reach the target with a very high many-body fidelity. For comparison, conventional QAOA appears insufficient to prepare the target state for the circuit depth of q = 18 (p = 9). The advantage of CD-QAOA is also visible in the entanglement entropy density curve: QAOA can easily lead to volumelaw entanglement, while CD-QAOA manages to generate as little entanglement as needed for the target state.

The discrepancy between conventional QAOA and CD-QAOA is best visible in the FM state preparation at  $\Delta/J=-2.0$ . In this case, a naïve application of QAOA with the set  $\mathcal{A}_{\text{QAOA}}=\{X|X+Y|Y,Z|Z\}$  is a priori doomed to fail: starting from the initial AFM state, which is orthogonal to the target FM manifold, the resulting QAOA unitaries leave the target AFM manifold invariant; in other words, transitions between the initial and the target states are forbidden by selection rules within the QAOA dynamics. Therefore, the many-body fidelity remains zero at all times during the QAOA evolu-

tion. The energy and entanglement entropy curves certify that the state does undergo nontrivial dynamics: similar to the XY state, QAOA creates volume-law entanglement and cannot reach the FM ground state manifold in energy, while CD-QAOA is well-behaved and sufficient to prepare the target. The CD-QAOA protocol sequence is shown in Table V(b): while we do not discern an obvious pattern, we emphasize that this time the RL agent makes use of both single-particle and two-body gauge potential terms.

Last, we show the system-size scaling of the energy curves for the three target states in Fig. 21(b-d). Similar to the spin-1/2 Ising chain, we find very little system-size dependence for the Haldane (b) and XY states (c). However, we cannot extrapolate the results to the thermodynamic limit due to the relatively small system sizes we were able to investigate. system-size effects are more pronounced for the ferromagnetic state (d), which is the one furthest away in Hilbert space from the initial perfect antiferromagnet.

Last, we mention in passing that we do not show results on preparing the AFM ground state at  $\Delta/J=2.0$  since this problem is somewhat trivial: indeed, starting from a perfect AFM in the z-direction, the AFM ground state of the spin-1 Heisenberg model can be easily reached even using adiabatic evolution because it lies within the AFM phase.

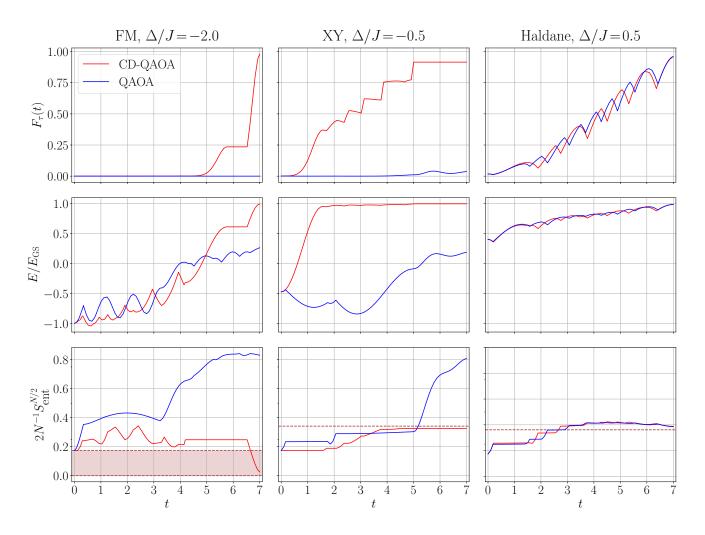


FIG. 20. Anisotropic Heisenberg spin-1 chain: time evolution generated by the protocol given by CD-QAOA (blue line), and conventional QAOA (red line) for the three target states, corresponding to the ferromagnetic  $(\Delta/J=-2.0)$ , XY  $(\Delta/J=-0.5)$ , and Haldane  $(\Delta/J=0.5)$  target state, respectively. Three quantities are shown: many-body fidelity (first row), energy ratio (second row), and the entanglement entropy density of the half chain (third row). The horizontal dashed line in the entanglement entropy curve shows the value in the target state, while the shaded area for the FM state denotes that in the span of the doubly degenerate ground state manifold. The protocols correspond to the duration T=7 in Fig. 3. The related CD-QAOA protocol sequences are given in Table V(b) [ferromagnetic  $(\Delta/J=-2.0)$ ], Table V(c) [XY  $(\Delta/J=-0.5)$ ] and Table V(d) [Haldane  $(\Delta/J=0.5)$ ]. The simulation parameters are the same as in Fig. 3.

#### 3. Lipkin-Meshkov-Glick model

In the main text, we also introduced the ferromagnetic Lipkin-Meshkov-Glick (LMG) model, described by the total spin Hamiltonian

$$H = -\frac{J}{N}(S^x)^2 + h\left(S^z + \frac{N}{2}\right).$$

Figure 22 shows the comparison between CD-QAOA and QAOA for two more values of h/J = 0.1 (deep in the ferromagnetic regime), and h/J = 0.9 (close to the critical point at h/J = 1.0). While the behavior for h/J = 0.1 is qualitatively similar to h/J = 0.5 (discussed in the main text), we do see that close to the critical point

the two-body gauge potential terms  $\hat{XY}$  and  $\hat{ZY}$  may offer some degree of improvement below the quantum speed limit, as compared to using only using the single-body  $\hat{Y}$  term. We mention in passing that we observed a stronger system-size dependence in the optimal protocol found by the RL agent in the immediate vicinity of the critical point  $h_c/J=1$ .

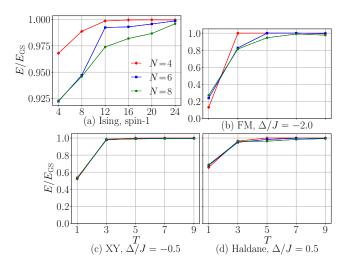


FIG. 21. system-size scaling of the energy minimization against protocol duration T for different system sizes N: (a) spin-1 Ising chain, (b-d) anisotropic Heisenberg spin-1 chain for  $\Delta/J=-2.0,\,\Delta/J=-0.5,\,\Delta/J=0.5$ , respectively. Note that the y-axis scale is different for the spin-1 Ising model in panel (a). The model parameters are the same as in (a) Fig. 7 and (b-d) Fig. 3 correspondingly.

# 4. Spin-1 Ising Chain

Finally, let us turn to the spin-1 Ising chain:

$$H(\lambda) = \lambda(t)H_1 + H_2,$$

$$H_1 = \sum_{j=1}^{N} J S_{j+1}^z S_j^z + h_x S_j^x,$$

$$H_2 = \sum_{j=1}^{N} h_z S_j^z,$$

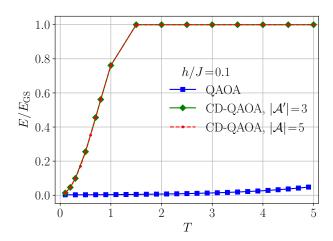
$$(F5)$$

see main text for discussion of the model parameters. Using this model, we compare four state preparation techniques: CD-QAOA, conventional QAOA, CD-driving using a variational gauge potential, and adiabatic evolution.

In order to compare these four methods, we first investigate their energy budget, i.e. the amount of energy required by the corresponding protocols. This is necessary, since variational CD-driving does not put any constraints on the magnitude of the expansion parameters  $\beta_j(\lambda)$  [cf. App. E], and we know that larger energies (i.e. generators of unitaries  $H_j$  with large norms) in general allow for a faster population transfer. To measure quantitatively the energy budget of a protocol, we use the average energy density along the protocol trajectory

$$\mathcal{N} = \frac{1}{T} \int_0^T dt \frac{\|H(t)\|}{N}, \tag{F6}$$

where H(t) is a unified notation for the continuous protocols in the case of adiabatic or CD driving, and the piecewise-constant (in time) sequences in CD-QAOA and



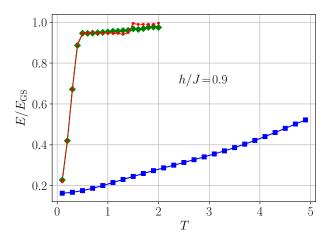


FIG. 22. LMG model: energy minimization against protocol duration T using conventional QAOA (blue square) and CD-QAOA (red dashed line, green solid line). The model parameters are the same from the settings in Fig. 5 but for h/J=0.1 (top panel), and h/J=0.9 (bottom panel).

conventional QAOA;  $\|H\|$  denotes the Hilbert-Schmidt norm of the operator H. Since we are interested in manybody systems, it is also natural to look at the energy density, i.e.  $\|H(t)\|/N$ . Figure 23 [bottom] shows that  $\mathcal N$  is on a similar scale for all four methods within the range of durations of interest, which allows for a meaningful comparison between them. As expected, CD-driving approaches adiabatic driving at large T, since the gauge potential term comes with a pre-factor  $\lambda$  which vanishes for  $T \to \infty$ ; in the opposite limit of  $T \to 0$ , the energy budget of CD-driving blows up, as a result of  $\beta_j(\lambda)$  being unconstrained.

In Fig. 23 [top], we see that the many-body fidelity, associated with the protocols obtained using energy density minimization, increases the performance contrast between the performance of the different methods [cf. Fig. 7, main text]. Since the fidelity is defined as the overlap square of the final with the target states [Eq. (F3)], like

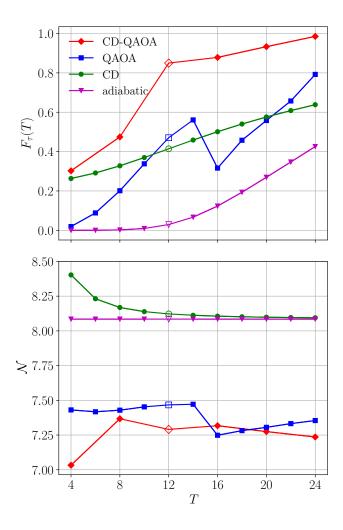


FIG. 23. Spin-1 Ising model: energy minimization against different protocol duration T for four different optimization methods: CD-QAOA (red line), conventional QAOA (blue line), variational gauge potential (green) and adiabatic evolution (magenta). Two associated quantities are shown: manybody fidelity  $F_{\tau}$  (top) and normalized time-averaged energy density  $\mathcal N$  over the protocol (bottom). The empty symbols mark the duration for which the evolution of physical quantities is shown in Fig. 24. The parameters are the same as in Fig. 7.

the entanglement entropy, it is insensitive to any specific observable; this implies that CD-QAOA outperforms the other three methods on all observables, not just energy. This is anticipated, because CD-QAOA combines the variational power of QAOA with physical insights from CD driving. Despite its better performance, notice how CD-QAOA also has a smaller energy budget than either of CD- and adiabatic driving.

To demonstrate the nonequilibrium character of the optimal protocols found by the RL agent in this setup, we fix T=12, and look at the time evolution of the energy, the fidelity, and the entanglement entropy within the learned protocol, cf. Fig. 24. While the protocol sequence [Table V(a)] appears impenetrable, we remark that (i) the RL agent makes use of both single-particle and twobody gauge potential terms, and (ii) some step durations  $\alpha_i$  are found to vanish identically, suggesting that the value of q may be reduced. As anticipated, the behavior of the dynamics generated by the CD and adiabatic driving is smooth, in contrast to the circuit-like piece-wise continuous curves of QAOA and CD-QAOA. The highly non-monotonic behavior of the energy curve shows that the CD-QAOA dynamics can be highly nonequilibrium: this likely stems from the RL objective [cf. App. A] – the total expected return: the agent only cares about maximizing the reward at t = T and is insensitive to any intermediate values. This allows the agent to drive the system through various states which are very far away from the target (e.g. w.r.t. the fidelity) [Curiously, these bad-energy states are all distinct, since they have different entanglement entropy, and the system does not visit the same quantum state twice during the evolution. The non-smooth and non-monotonic behavior of the CD-QAOA solution raises the question about how robust the protocol is, to small external perturbations – a topic of future studies.

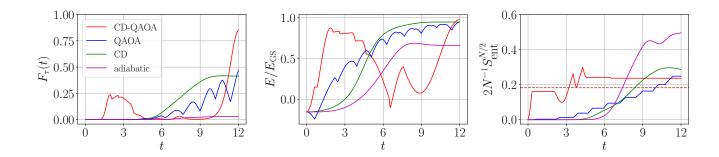


FIG. 24. Spin-1 Ising model: time evolution generated by the four different methods: CD-QAOA (red), conventional QAOA (blue), CD driving using the variational gauge potential (green) and adiabatic evolution (magenta). The three quantities are shown: the many-body fidelity (left), energy (middle), and entanglement entropy of the half chain (right). The protocols correspond to the empty symbols during T=12 in Fig. 7. We compare The horizontal dashed line in the entanglement entropy curve shows the value in the target state. The CD-QAOA protocol sequence is given in Table V(a). The model parameters are the same as in Fig. 7.

# (b) Ferromagnetic $(\Delta/J = -2.0)$

( ) 0 1			
Hamiltonian	Duration		
X Y	0.312		
$\dot{Y}$	0.299		
Z	0.216		
Y	0.717		
Z	0.000		
Y	0.537		
Z Z+X	0.477		
Y	0.054		
Z Z+X	0.657		
Z	0.000		
Z Z+X	0.269		
Y Z	0.274		
Z Z+X	0.478		
Y Z	0.372		
Z Z+X	0.000		
Z	1.794		
X Y	0.072		
Z	0.039		
Y	1.007		
Z	4.426		

(a) Ising spin-1

short-hand notation	Duration
Y Z-YZ	0.122
X X+Y Y	0.178
YZ	0.027
Z Z	0.376
Y Z-YZ	0.234
Z Z	0.000
X X+Y Y	0.323
Z Z	0.284
Y Z-YZ	0.366
Z Z	0.000
X X+Y Y	0.314
Z Z	0.188
Y Z-YZ	0.535
Y	0.001
X X	0.342
Z Z	0.105
Y Z-YZ	0.538
X X	0.208
Y	0.000
Z Z	0.051
Y	0.658
Y Z-YZ	0.002
Y	0.900
$Z \ Y$	0.771
	0.005
X Y-XY	0.474
Y Z-YZ	0.000
X X+Y Y	0.000

1	(c)	$\mathbf{X}\mathbf{Y}$	(Λ	/ 1		_0	۲,
(	$^{\rm c}$	AI	$\Delta$	/J	=	−υ.	Э

short-hand notation	Duration
Y	0.795
X X	0.000
Y	0.772
X X+Y Y	0.143
X X	0.383
Y	0.001
X X+Y Y	0.284
X X	0.180
X X+Y Y	0.467
X X	0.113
X X+Y Y	0.635
X X	0.097
X X+Y Y	0.617
Y	0.000
Z Z	0.162
X X+Y Y	0.265
X X	0.092
Z	1.995

(d) Haldane ( $\Delta/J=0.5$ )

short-hand notation	Duration
X X+Y Y	0.149
X X	0.000
X X+Y Y	0.052
Z Z	1.376
X X+Y Y	0.313
Z Z	0.668
X X+Y Y	0.187
Z Z	0.723
X X+Y Y	0.289
Z Z	0.528
X X+Y Y	0.218
Z Z	0.561
X X+Y Y	0.254
Z Z	0.684
X X+Y Y	0.360
Z Z	0.639
X X	0.000
Z	0.000

TABLE V. Ising spin-1 chain and Anisotropic Heisenberg spin-1 chain: the protocol sequences and corresponding durations given by CD-QAOA. The protocol (a) correspond to Ising spin-1 in Fig. 24; the (b), (c), (d) three sequences correspond to the three phases in the same setting as Fig. 20. The short-hand notation is the same in Table I. We use a shaded cell background whenever terms from the CD gauge potential are used in the protocol sequence. Terms of zero durations are marked in light grey.