

# Exploring fixed points and eigenstates of quantum systems with reinforcement learning

María Laura Olivera-Atencio,<sup>1,\*</sup> Jesús Casado-Pascual,<sup>1,2,†</sup> and Denis Lacroix<sup>3,‡</sup>

<sup>1</sup>*Física Teórica, Universidad de Sevilla, Apartado 1065, E-41080 Sevilla, Spain*

<sup>2</sup>*Multidisciplinary Unit for Energy Science, Universidad de Sevilla, E-41080 Sevilla, Spain*

<sup>3</sup>*Université Paris-Saclay, CNRS/IN2P3, IJCLab, 91405 Orsay, France*

(Dated: November 24, 2025)

We introduce a reinforcement learning algorithm designed to identify the fixed points of a given quantum operation. The method iteratively constructs the unitary transformation that maps the computational basis onto the basis of fixed points through a reward–penalty scheme based on quantum measurements. In cases where the operation corresponds to a Hamiltonian evolution, this task reduces to determining the Hamiltonian eigenstates. The algorithm is first benchmarked on random Hamiltonians acting on two and three qubits and then applied to many-body systems of up to six qubits, including the transverse-field Ising model and the all-to-all pairing Hamiltonian. In both cases, the algorithm is demonstrated to perform successfully; in the pairing model, it can also reveal hidden symmetries, which can be exploited to restrict learning to specific symmetry sectors. Finally, we discuss the possibility of post-selecting high-fidelity states even when full convergence has not been reached.

Keywords: quantum reinforcement learning, quantum computing, quantum algorithms

## I. INTRODUCTION

Quantum machine learning (QML) has emerged as a fertile ground for research in recent years, exploring the interplay between quantum mechanics and machine learning [1–6]. This field aims to leverage quantum technologies to enhance various machine learning tasks, including classification problems [7], quantum control [8], or enhancing of quantum synchronization processes [9]. Diverse quantum algorithms have been proposed and, in some cases, experimentally implemented, demonstrating potential speedups over classical approaches [10, 11].

Reinforcement learning (RL) represents a fundamental approach within the field of machine learning [12]. In contrast to supervised and unsupervised learning, which mainly focus on identifying patterns in data [13–17], RL addresses complex problems by framing them as the search for an effective sequence of decisions that maximizes long-term rewards. This is accomplished through the dynamic interaction between an agent and its environment. The learning process is driven by repeated interactions, where actions are either rewarded or penalized based on a predefined policy, aiming to maximize cumulative rewards [12, 18, 19]. The intersection between this iterative trial-and-error methodology and quantum mechanics has given rise to a diverse body of work ranging from fundamental studies [20], studies focused on the use of quantum properties to enhance learning tasks [21, 22], to specialized tasks such as preparation of quantum states or the search for quantum error correction strategies [23–25].

In recent years, this paradigm has continued to inspire further developments in quantum RL [26]. For instance, high fidelity state control has been achieved using quantum RL for optimization when full observation states are infeasible [27]. Policy iteration methods have been explored to improve the performance of quantum reinforcement algorithms in the context of infinite-horizon discounted problems [28]. Furthermore, challenges such as managing continuous action spaces [29], as well as the implementation of a quantum decision problem in a quantum computer, have been successfully addressed [30], including tackling eigenvalue problems in the presence of dissipation [31–33].

The present work introduces an RL algorithm designed to approximately learn the fixed points of a given quantum operation. More precisely, the proposed algorithm identifies the unitary transformation that maps the computational basis onto the basis composed of the pure states that remain invariant under the action of a specified quantum operation. When the operation corresponds to unitary dynamics generated by a Hamiltonian, the fixed points coincide with its eigenstates, and the algorithm provides the associated eigenbasis. This approach is inspired by previous work [34], which focused on identifying the eigenvectors of unknown Hamiltonians. A key difference is that our algorithm learns all fixed points simultaneously, whereas in Ref. [34] the eigenstates were determined sequentially, one at a time. Instead of identifying individual eigenstates, our method learns the global unitary transformation that connects the computational basis to the fixed-point basis, providing a more compact and parallelizable framework. Moreover, while the previous approach was restricted to unitary processes, the method presented here is, in principle, applicable to general quantum operations, including non-unitary evolutions. Despite this broader applicability, in the present

\* molivera1@us.es

† jcasado@us.es

‡ lacroix@ijclab.in2p3.fr

work we focus exclusively on the unitary case, leaving the analysis of noise and dissipative effects for future investigations. Preliminary results in this direction can be found in Refs. [31–33].

It is worth noting that the problem of determining Hamiltonian eigenstates has been addressed through a wide variety of approaches. Specifically, in recent years, significant efforts have been devoted to the development of hybrid quantum-classical variational techniques aimed at constructing approximate ground states of complex many-body quantum systems [35–37]. Variational methods based on parametrized trial wavefunctions have been extensively explored in different areas of physics and chemistry [38, 39], offering the advantage, particularly for near-term applications, of being easier to implement on imperfect quantum devices. They have been successfully applied to many-body systems of moderate size; however, as the number of qubits increases, the optimization task typically becomes more challenging and may fail due to the complexity of the energy landscape, with the onset of the so-called “barren plateau” phenomenon being a paradigmatic example [40]. Extensive efforts have been made to extend variational approaches and, in some cases, complement them with quantum subspace diagonalization techniques to also provide accurate descriptions of excited states [41–73]. Our aim with the method presented here is not to compete with existing variational algorithms, but rather to explore a non-variational alternative that could be useful for addressing such problems in many-body systems. For this reason, in this paper we also apply our RL algorithm to two types of interacting-particle models that are relevant, respectively, for the condensed matter and nuclear physics communities.

The remainder of the paper is organized as follows. In Sec. II, we provide a detailed description of the algorithm. We first define the problem statement and goal, specifying the task of learning the fixed points of a quantum operation. We then present the RL protocol step by step, detailing the iterative procedure and the reward–punishment strategy. Next, we discuss strategies to prevent convergence to spurious fixed points and describe a fine-tuning procedure after initial convergence to further improve fidelities. Metrics for assessing convergence and accuracy are also introduced, followed by initial benchmarks on random Hamiltonians that illustrate the performance of the method. Section III focuses on applications to physical systems, namely the transverse-field Ising model (TFIM) and the pairing model, and discusses extensions such as symmetry-restricted RL and post-selection of high-fidelity states. Finally, Sec. IV summarizes our main findings and outlines possible directions for future work.

## II. ALGORITHM DESCRIPTION

### A. Problem statement and goal

Consider a quantum system with a state space of dimension  $d$  and computational basis  $\mathbf{C} = \{|j\rangle\}_{j \in \{0, \dots, d-1\}}$ . Let  $\mathbf{B} = \{|\Phi_\alpha\rangle\}_{\alpha \in \{0, \dots, d-1\}}$  be another, unknown orthonormal target basis that also spans the state space of the system. We assume that the only information available about the basis  $\mathbf{B}$  is that it is left invariant by a quantum operation  $\mathcal{E}$ , meaning that  $\mathcal{E}(|\Phi_\alpha\rangle\langle\Phi_\alpha|) = |\Phi_\alpha\rangle\langle\Phi_\alpha|$  for all  $\alpha \in \{0, \dots, d-1\}$ . With this hypothesis our objective is to construct, at least approximately, a unitary transformation that maps the computational basis  $\mathbf{C}$  onto the target basis  $\mathbf{B}$  (up to a permutation of their elements).

A particular case—on which we will focus later in this work—is when the target basis  $\mathbf{B}$  consists of the (unknown) stationary states of a certain Hamiltonian  $H$ , that is,  $H|\Phi_\alpha\rangle = E_\alpha|\Phi_\alpha\rangle$  for all  $\alpha \in \{0, \dots, d-1\}$ , where  $E_\alpha$  is the eigenenergy associated with the stationary state  $|\Phi_\alpha\rangle$ . In this setting, the target basis  $\mathbf{B}$  remains invariant under the action of the quantum operation

$$\mathcal{E}(\bullet) = U(\tau) \bullet U^\dagger(\tau), \quad (1)$$

where  $U(\tau) = e^{-i\tau H/\hbar}$  is the unitary time-evolution operator from time 0 to  $\tau$ , with  $\tau$  being a time interval whose value is, for now, arbitrary. As we will see later, the particular choice of the parameter  $\tau$  plays a crucial role in the correct functioning of the algorithm.

The goal of our quantum RL method is to iteratively construct a sequence of unitary matrices  $\{D_k\}$ , where  $k$  labels the iteration number, such that

$$\lim_{k \rightarrow +\infty} D_k |j\rangle \approx |\Phi_\alpha\rangle \quad (2)$$

for all  $|j\rangle$ , with  $|\Phi_\alpha\rangle$  representing some state of the basis  $\mathbf{B}$ . In other words, in the limit  $k \rightarrow +\infty$ ,  $D_k$  approximately maps the basis  $\mathbf{C}$  onto the basis  $\mathbf{B}$ , up to a certain permutation of its vectors.

### B. Technical details on the iterative process

The unitary transformation generated in the  $(k+1)$ th iteration,  $D_{k+1}$ , is obtained from the unitary operator  $D_k$ , corresponding to iteration  $k$ , according to the update rule

$$D_{k+1} = D_k \prod_{j=0}^{d-2} \prod_{l=j+1}^{d-1} D_k^{(j,l)}, \quad (3)$$

where the  $d(d-1)/2$  operators  $D_k^{(j,l)}$  are unitary transformations. The index  $(j,l)$  refers to the  $\mathbf{C}$  basis states. In this product, the operators are arranged from left to right according to the increasing sequence of index pairs  $(j,l)$ , ordered lexicographically. Specifically,  $(j_1, l_1)$  precedes  $(j_2, l_2)$  if either  $j_1 < j_2$ , or  $j_1 = j_2$  and  $l_1 < l_2$ .

We assume that the unitary operators  $D_k^{(j,l)}$  are taken as rotations acting on the two-dimensional subspace spanned by  $|j\rangle$  and  $|l\rangle$ . These rotations can be expressed in terms of the Pauli-type operators

$$\begin{cases} X^{(j,l)} = |j\rangle\langle l| + |l\rangle\langle j|, \\ Y^{(j,l)} = -i(|j\rangle\langle l| - |l\rangle\langle j|), \\ Z^{(j,l)} = |j\rangle\langle j| - |l\rangle\langle l|, \end{cases} \quad (4)$$

which serve as their generators. Specifically, the operators  $D_k^{(j,l)}$  are expressed as

$$D_k^{(j,l)} = e^{-i\beta_k^{(j,l)} Y^{(j,l)}/2} e^{-i\gamma_k^{(j,l)} Z^{(j,l)}/2} e^{-i\alpha_k^{(j,l)} X^{(j,l)}/2}, \quad (5)$$

where the angles  $\alpha_k^{(j,l)}$ ,  $\beta_k^{(j,l)}$ , and  $\gamma_k^{(j,l)}$  are chosen depending on whether the corresponding action is rewarded or penalized. In the case of a reward, all angles are set to zero, so that  $D_k^{(j,l)}$  reduces to the identity operator  $I$ . This corresponds to no exploration, i.e., pure exploitation in RL terms. Conversely, in the case of a penalty, the angles are drawn pseudo-randomly and uniformly within intervals that are dynamically adjusted during the learning process, allowing exploration of the space of unitary operators.

Schematically, these intervals are expanded in response to penalties and contracted in response to rewards, thereby regulating the extent of exploration in the space of unitary operators. The degree of exploration at iteration  $k$  is controlled via a set of  $d(d-1)/2$  exploration parameters  $w_k^{(j,l)} \leq 1$ , whose evolution with the number of iterations depends on the cumulative number of rewards and punishments received up to that point (see the discussion below for further details).

### C. RL algorithm step-by-step

The algorithm is implemented on a set of  $d$   $d$ -dimensional systems (qudits), each initialized in a distinct state of the computational basis. Specifically, the initial density operator of the full system is

$$\rho_1 = |0\rangle\langle 0| \otimes |1\rangle\langle 1| \otimes \cdots \otimes |d-1\rangle\langle d-1|. \quad (6)$$

The iterative process begins with the initial assignments  $D_1 = I$  and  $w_1^{(j,l)} = 1$  for all index pairs  $(j,l)$  such that  $0 \leq j < l \leq d-1$ . At each iteration, the values of  $D_{k+1}$  and the parameters  $w_{k+1}^{(j,l)}$  are updated based on those from the previous iteration,  $D_k$  and  $w_k^{(j,l)}$ , following the procedure described below:

1. **Apply  $D_k$ :** The unitary transformation  $D_k$  is applied to each qudit individually, yielding the density operator

$$\rho_k = \rho_k^{(0)} \otimes \rho_k^{(1)} \otimes \cdots \otimes \rho_k^{(d-1)}, \quad (7)$$

with  $\rho_k^{(j)} = D_k |j\rangle\langle j| D_k^\dagger$ .

2. **Apply  $\mathcal{E}$ :** Next, the quantum operation  $\mathcal{E}$  is applied to each qudit individually, resulting in the updated density operator

$$\tilde{\rho}_k = \tilde{\rho}_k^{(0)} \otimes \tilde{\rho}_k^{(1)} \otimes \cdots \otimes \tilde{\rho}_k^{(d-1)}, \quad (8)$$

with  $\tilde{\rho}_k^{(j)} = \mathcal{E}(\rho_k^{(j)})$ .

3. **Undo  $D_k$ :** The next step consists in reversing the unitary transformation applied in step 2 for each qudit. This yields the density operator

$$\tilde{\tilde{\rho}}_k = \tilde{\tilde{\rho}}_k^{(0)} \otimes \tilde{\tilde{\rho}}_k^{(1)} \otimes \cdots \otimes \tilde{\tilde{\rho}}_k^{(d-1)}, \quad (9)$$

with  $\tilde{\tilde{\rho}}_k^{(j)} = D_k^\dagger \tilde{\rho}_k^{(j)} D_k$ .

4. **Measure:** Then, a projective measurement is performed on each qudit in the computational basis, with the result for the  $j$ th qudit denoted by  $m_k^{(j)} \in \{0, 1, \dots, d-1\}$ . It is worth noting that if, at a given iteration  $k$ , the target were exactly reached—namely, if  $D_k$  mapped the computational basis  $\mathbf{C}$  onto the basis  $\mathbf{B}$  (up to a permutation of their elements)—then the density operator  $\rho_k$  would remain invariant under the quantum operation  $\mathcal{E}$ . In that case, we would have  $\tilde{\rho}_k = \rho_k$ , and upon reversing the transformation  $D_k$ , the initial density operator would be recovered, i.e.,  $\tilde{\tilde{\rho}}_k = \rho_1$ . If this condition held, the measurement outcomes would satisfy the equality  $m_k^{(j)} = j$  for all  $j \in \{0, \dots, d-1\}$  with probability one. This is equivalent to stating that, for each label  $j$  and any  $l \neq j$ , the condition  $m_k^{(j)} \neq l$  would hold with probability one. Any outcome consistent with this condition would indicate that the goal may have been achieved and should therefore be rewarded, whereas any violation would imply that the objective has not yet been reached and should be penalized.

5. **Determination of  $D_k^{(j,l)}$  and  $w_{k+1}^{(j,l)}$  (reward/punishment):** Depending on the outcomes of the previous measurements, the following procedure is carried out for each pair of qudits  $(j,l)$  with  $0 \leq j < l \leq d-1$ :

i) **Both correct (double reward):** If  $m_k^{(j)} \neq l$  and  $m_k^{(l)} \neq j$ , the condition is satisfied for both qudits. In this case, a double reward is applied by reducing the corresponding exploration parameter according to the rule

$$w_{k+1}^{(j,l)} = r^2 w_k^{(j,l)}, \quad (10)$$

where  $r \in (0, 1)$  is a parameter referred to as the reward rate. Furthermore, the three angles in Eq. (5) are set to zero, so that  $D_k^{(j,l)}$  reduces to the identity operator, i.e.,  $D_k^{(j,l)} = I$ .

ii) **One correct and one incorrect (reward & punishment):** If  $m_k^{(j)} \neq l$  and  $m_k^{(l)} = j$ , or  $m_k^{(j)} = l$  and  $m_k^{(l)} \neq j$ , then one qudit satisfies the condition while the other violates it. In this case, a mixed reward–punishment update is performed. Specifically, the

exploration parameter is updated according to the rule

$$w_{k+1}^{(j,l)} = \min(rp w_k^{(j,l)}, 1), \quad (11)$$

where  $p > 1$  is a parameter referred to as the punishment rate. Thus, the exploration parameter decreases if  $rp < 1$ , remains unchanged if  $rp = 1$ , and increases if  $rp > 1$ , with the min function preventing it from exceeding 1. The unitary operator  $D_k^{(j,l)}$  is then constructed as in Eq. (5), with the angles  $\alpha_k^{(j,l)}$ ,  $\beta_k^{(j,l)}$ , and  $\gamma_k^{(j,l)}$  drawn pseudo-randomly and uniformly within the interval  $[-\pi w_k^{(j,l)}, \pi w_k^{(j,l)}]$ .

iii) **Both incorrect (double punishment):** If  $m_k^{(j)} = l$  and  $m_k^{(l)} = j$ , then both qudits violate the condition. In this case, a double punishment is applied, with the exploration parameter updated according to the rule

$$w_{k+1}^{(j,l)} = \min(p^2 w_k^{(j,l)}, 1). \quad (12)$$

The construction of  $D_k^{(j,l)}$  follows the same procedure as in the case ii).

**6. Determination of  $D_{k+1}$ :** Once the above procedure has been applied to compute the operators  $D_k^{(j,l)}$  for all qudit pairs, the operator  $D_{k+1}$  is obtained from  $D_k$  using Eq. (3).

**7. Reset:** If the measurement outcome  $m_k^{(j)}$  differs from  $j$ , the unitary  $|m_k^{(j)}\rangle\langle j| + |j\rangle\langle m_k^{(j)}| + \sum_{l \neq j, m_k^{(j)}} |l\rangle\langle l|$  is applied to the corresponding qudit to reset it to its initial computational basis state. Once all necessary qudits have been reinitialized in this way, the algorithm resumes from step 1 with the iteration index updated to  $k+1$ .

The algorithm described above is said to converge if all exploration parameters  $w_k^{(j,l)}$  approach zero as the number of iterations  $k$  becomes sufficiently large. In this case, all the operators  $D_k^{(j,l)}$  tend to the identity, and the unitary transformation  $D_k$  approaches a fixed value. To determine in practice whether the algorithm has converged, it is useful to introduce a sufficiently small threshold value for the exploration parameter, denoted by  $w_{\text{th}}$ . Additionally, we define the maximum exploration parameter at iteration  $k$  as

$$w_k^{(M)} = \max_{0 \leq j < l \leq N-1} w_k^{(j,l)}. \quad (13)$$

We say that the algorithm has converged if  $w_k^{(M)} < w_{\text{th}}$  after a certain number of iterations  $k$ . A schematic representation of the algorithm described above is shown in Fig. 1.

#### D. Fixed-point ambiguities: ensuring convergence to the target basis

As noted in Sec. IIA, the purpose of the algorithm is to construct a unitary transformation that maps the

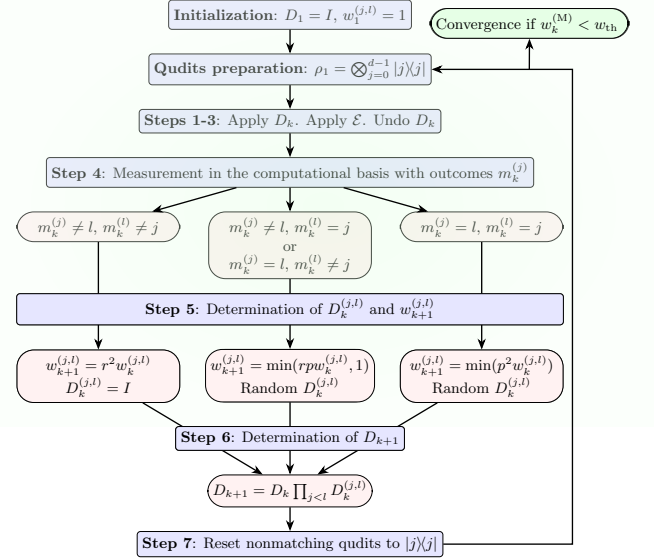


FIG. 1. Schematic representation of the proposed RL algorithm.

computational basis  $\mathcal{C}$  onto the target basis  $\mathcal{B}$ , under the assumption that  $\mathcal{B}$  consists of pure states left invariant by a certain quantum operation  $\mathcal{E}$ . If, however,  $\mathcal{B}$  is not the only invariant set of pure states under such an operation, the algorithm may instead converge to a unitary mapping  $\mathcal{C}$  onto a different invariant set, rather than the intended target basis.

To illustrate this point, let us consider the case where the target basis is formed by the stationary states of a Hamiltonian  $H$ , and the quantum operation  $\mathcal{E}$  is given by Eq. (1). In this situation, it is straightforward to verify that, for certain values of the evolution time  $\tau$ , pure states other than the stationary states of  $H$  may also remain invariant under  $\mathcal{E}$ . A necessary and sufficient condition for a pure state  $|\Psi\rangle\langle\Psi|$  to be invariant is that the survival probability after evolution for a time  $\tau$ , defined as

$$P = \langle\Psi|U(\tau)|\Psi\rangle\langle\Psi|U^\dagger(\tau)|\Psi\rangle = |\langle\Psi|U(\tau)|\Psi\rangle|^2, \quad (14)$$

equals unity. Expanding  $|\Psi\rangle$  in the eigenbasis of  $H$  as  $|\Psi\rangle = \sum_{\alpha=0}^{d-1} c_\alpha |\Phi_\alpha\rangle$ , one finds

$$P = 1 - 4 \sum_{\alpha=0}^{d-2} \sum_{\beta=j+1}^{d-1} |c_\alpha|^2 |c_\beta|^2 \left[ \sin\left(\frac{\omega_{\alpha,\beta}\tau}{2}\right) \right]^2, \quad (15)$$

with  $\omega_{\alpha,\beta} = (E_\alpha - E_\beta)/\hbar$ . According to this expression, if  $\tau$  is such that  $\omega_{\alpha,\beta}\tau$  is an integer multiple of  $2\pi$ , then  $\sin(\omega_{\alpha,\beta}\tau/2) = 0$  and any superposition of the form  $|\Psi\rangle = c_\alpha |\Phi_\alpha\rangle + c_\beta |\Phi_\beta\rangle$  remains invariant under  $\mathcal{E}$ . However, unless  $\omega_{\alpha,\beta} = 0$  (degenerate eigenvalues), such states are not true eigenstates of  $H$  when both coefficients  $c_\alpha$  and  $c_\beta$  are nonzero.

This situation leads to an ambiguity in the learning dynamics: For certain values of  $\tau$ , the algorithm cannot

reliably distinguish these superpositions from the true eigenstates based solely on their invariance under the quantum operation in Eq. (1). Moreover, the issue is not limited to cases where  $\omega_{\alpha,\beta}\tau$  is exactly resonant, i.e., an integer multiple of  $2\pi$ ; even when  $\omega_{\alpha,\beta}\tau$  takes near-resonant values, the algorithm may struggle to differentiate between genuine eigenstates and spurious invariant superpositions, potentially leading to convergence toward an incorrect solution.

Since we do not know *a priori* whether a fixed  $\tau$  satisfies this problematic condition, a more robust strategy is to avoid using the same  $\tau$  in every iteration. Instead, a different value of  $\tau$  is randomly selected at each step from a suitably chosen interval. In this way, even if certain  $\tau$  values prevent the algorithm from distinguishing some superpositions from true eigenstates, other values will not, thereby mitigating the occurrence of spurious eigenstates in the averaged behavior over multiple iterations. For each pair of stationary states with distinct energies (i.e.,  $\omega_{\alpha,\beta} \neq 0$ ), the sinusoidal terms in Eq. (15) have a period  $2\pi/\omega_{\alpha,\beta}$ , so the longest period among them is  $\tau_{\max} = 2\pi/\omega_{\min}$ , where  $\omega_{\min}$  denotes the smallest nonzero  $\omega_{\alpha,\beta}$ . Sampling  $\tau$  from an interval of width on the order of  $\tau_{\max}$ —for instance,  $\tau \in [0, \tau_{\max}]$ —helps to avoid systematic reinforcement of spurious solutions across iterations. If the spectrum of  $H$  is unknown, as is the case in most practical situations,  $\tau_{\max}$  can be roughly estimated from the characteristic energy scales of the problem.

### E. Fine-tuning the algorithm after initial convergence

A way to improve the precision of the algorithm is to introduce a fine-tuning stage once an initial convergence has been reached. Specifically, when the algorithm is considered to have converged—i.e., when the maximum exploration parameter satisfies  $w_k^{(M)} < w_{\text{th}}$ —all exploration parameters  $w_k^{(j,l)}$  are reset to a new value  $w_r$ , chosen such that  $w_{\text{th}} < w_r < 1$ . This reinitialization introduces a new phase of exploration while preserving most of the knowledge already acquired.

The value of  $w_r$  should be carefully selected: If it is too close to 1, small errors in the measurements could undo the learning already achieved. On the other hand, if it is too close to the threshold  $w_{\text{th}}$ , the exploration will be too limited to allow any significant refinement. The optimal value of  $w_r$  must be therefore adjusted depending on the specific problem at hand.

If the reset procedure were applied indefinitely, the algorithm would never reach full convergence. To avoid this, the reset value  $w_r$  is gradually decreased after a fixed number of iterations. Specifically, once the iteration index exceeds a predetermined value  $k_0$ ,  $w_r$  is reduced linearly according to the rule  $w_r(k_M - k)/(k_M - k_0)$ , where  $k_M$  denotes the maximum number of iterations allowed.

This schedule enables finer exploration after the initial

convergence phase, giving the algorithm the flexibility to improve the solution while preserving the knowledge acquired in earlier stages.

### F. Assessing algorithm convergence and accuracy

Two key aspects must be considered when assessing the performance of the algorithm described above. First, the number of iterations required to achieve convergence should be as small as possible, since faster convergence reflects better algorithmic efficiency. Second, the algorithm must reproduce the target basis states  $\mathbf{B}$  with high accuracy, ideally within the minimum number of iterations.

As discussed earlier, convergence can be quantified through the dependence of the maximum exploration parameter  $w_k^{(M)}$  on the iteration index  $k$ . Regarding accuracy, in all the examples considered below the states of  $\mathbf{B}$  can be computed independently, without relying on the algorithm. This makes it possible to evaluate accuracy directly by calculating the overlaps between the algorithmic states  $\{D_k|j\rangle\}_{j \in \{0, \dots, d-1\}}$  and the independently obtained vectors of  $\mathbf{B}$ . Since it is not known *a priori* which vector of the target basis  $\mathbf{B}$  a given computational basis state  $|j\rangle$  will be mapped to, it is necessary to adopt a flexible accuracy measure. A suitable measure of the algorithm's accuracy is provided by the maximum square-root fidelities, defined for each  $j \in \{0, \dots, d-1\}$  as

$$f_k^{(j)} = \max_{\alpha \in \{0, \dots, d-1\}} |\langle \Phi_\alpha | D_k | j \rangle|, \quad (16)$$

which reflects the best overlap achieved by the transformed state  $D_k|j\rangle$  with any state of the target basis  $\mathbf{B}$  at iteration  $k$ .

It is important to note that the algorithm proposed here explores the space of unitary operators in a stochastic manner. This randomness is twofold: on one hand, the measurement outcomes in step 4 are intrinsically probabilistic; on the other hand, the angles  $\alpha_k^{(j,l)}$ ,  $\beta_k^{(j,l)}$ , and  $\gamma_k^{(j,l)}$  are generally chosen pseudo-randomly. Consequently, repeated applications of the algorithm generally produce different instances of  $D_k$ , each with varying accuracy. Even when the accuracy is high, in different realizations the same computational basis vector  $|j\rangle$  may be mapped to different vectors in the target basis  $\mathbf{B}$ . Therefore, to properly assess the algorithm's performance, it is convenient to consider a sufficiently large set of  $N_r$  realizations. For each iteration  $k$ , the arithmetic means of the maximum exploration parameter  $w_k^{(M)}$  and the maximum square-root fidelity  $f_k^{(j)}$  across realizations are computed and denoted as  $W_k^{(M)}$  and  $F_k^{(j)}$ , respectively. These averaged quantities provide a robust way to quantify both the convergence and accuracy of the algorithm over multiple stochastic realizations.

In what follows, we test the convergence and accuracy of the proposed algorithm by applying it to the calculation of the eigenstates of several representative Hamiltonian models. For this purpose, all Hamiltonians are rescaled to a dimensionless form,

$$\tilde{H} = \frac{H - E_{\min}I}{E_{\max} - E_{\min}}, \quad (17)$$

where  $E_{\min}$  and  $E_{\max}$  are the minimum and maximum eigenenergies of  $H$ , respectively. This guarantees that the spectrum of  $\tilde{H}$  lies between 0 and 1, with the lowest eigenvalue equal to 0 and the highest equal to 1. In addition, we define a dimensionless evolution time as  $\tilde{\tau} = \tau(E_{\max} - E_{\min})/\hbar$ , so that the time parameter naturally adapts to the spectral range of the Hamiltonian. To avoid convergence of the algorithm to spurious eigenstates, instead of keeping  $\tau$  (or its dimensionless version  $\tilde{\tau}$ ) fixed, at each iteration we randomly select its value within a suitable interval, as described in Sec. IID.

### G. Testing the algorithm on random Hamiltonians

To provide an initial illustration of the proposed algorithm, we apply it to the calculation of the eigenvectors of random Hamiltonians. In each realization, a different Hermitian matrix with randomly chosen complex entries is generated, and the associated Hamiltonian is rescaled according to Eq. (17). The algorithm is then executed for a sufficient number of iterations to ensure convergence, and the obtained states are compared with the exact eigenvectors computed by direct diagonalization, using the fidelities defined in Eq. (16). Repeating this procedure over  $N_r$  realizations finally allows us to compute the mean fidelities  $F_k^{(j)}$  and the mean maximum exploration parameter  $W_k^{(M)}$ .

Figure 2 shows the results obtained for a two-qubit system ( $d = 4$ ). Panels (a) and (b) display the behavior of  $F_k^{(j)}$  and  $W_k^{(M)}$  as a function of the iteration number  $k$  without the reset mechanism described in Sec. IIE, while panels (c) and (d) show the corresponding results when the reset is applied. In both cases, in each iteration the dimensionless evolution time  $\tilde{\tau}$  was uniformly sampled from the interval  $[0, 100]$  to avoid the issue discussed in Sec. IID. The remaining parameters are set to  $r = 0.9$ ,  $p = 2/r$ ,  $w_{\text{th}} = 0.005$ ,  $w_r = 0.01$ , and  $N_r = 100$ . The fidelities corresponding to the four possible states of the computational basis are represented in different colors. However, since the curves largely overlap, the colors cannot be easily distinguished.

As shown in Fig. 2, the algorithm attains high fidelities both with and without the reset mechanism, with a slight improvement when resets are applied. After convergence, the maximum and minimum fidelities across the computational basis states are nearly identical: without resets  $F_{\max} \approx F_{\min} \approx 0.97$ , whereas with resets  $F_{\max} \approx F_{\min} \approx 0.98$ . The average number of iterations

required for convergence is lower in the absence of resets, since the reset scheme necessarily introduces additional stages of exploration and, consequently, a higher number of iterations.

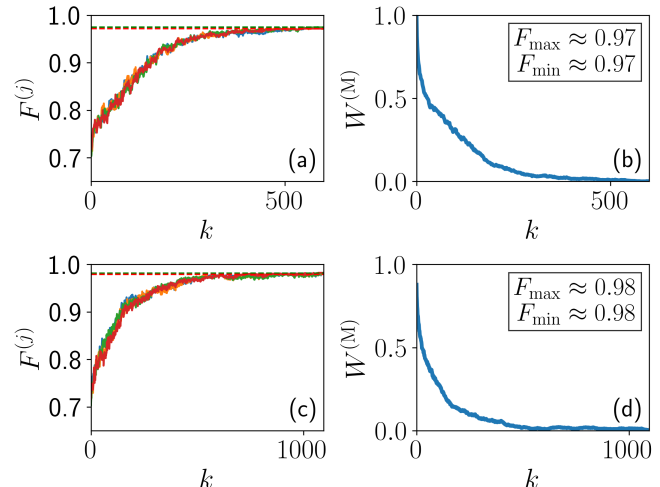


FIG. 2. Algorithm results for random two-qubit Hamiltonians ( $d = 4$ ). Panels (a) and (b) show the fidelities  $F_k^{(j)}$  and the maximum exploration parameter  $W_k^{(M)}$  as functions of the iteration number  $k$  without the reset mechanism described in Sec. IIE, while panels (c) and (d) show the results with reset. Each realization uses a different random Hamiltonian, and in each iteration the dimensionless evolution time  $\tilde{\tau}$  is uniformly sampled from  $[0, 100]$ . Parameters are  $r = 0.9$ ,  $p = 2/r$ ,  $w_{\text{th}} = 0.005$ ,  $w_r = 0.01$ , and  $N_r = 100$ . Fidelities for the four computational basis states are plotted in different colors; however, because the curves largely overlap, the colors are difficult to distinguish. In the left panels, dashed horizontal lines indicate the maximum and minimum fidelities,  $F_{\max}$  and  $F_{\min}$ .

In Fig. 3, we present results analogous to those in Fig. 2, but for a three-qubit system ( $d = 8$ ). The reset mechanism described in Sec. IIE substantially improves the fidelities, increasing them from  $F_{\min} \approx 0.90$  and  $F_{\max} \approx 0.91$  without resets to  $F_{\min} \approx 0.97$  and  $F_{\max} \approx 0.98$  with resets. While this improvement comes with an increase in the number of iterations required for convergence, the reset scheme consistently guides the algorithm toward high-quality results across all basis states, demonstrating its effectiveness in larger systems.

Random Hamiltonians of the type considered in this section have also been studied in the literature using a reinforcement learning algorithm different from the one presented here (see Ref. [34]). Providing a detailed description of the protocol of Ref. [34] is beyond the scope of this work. Nevertheless, although both approaches use measurement outcomes as the basis for reward and punishment, there are fundamental differences between the two reinforcement learning strategies. In particular, while the method of Ref. [34] searches for eigenstates sequentially, starting from an arbitrarily chosen one, our approach is designed to learn all eigenstates simultane-

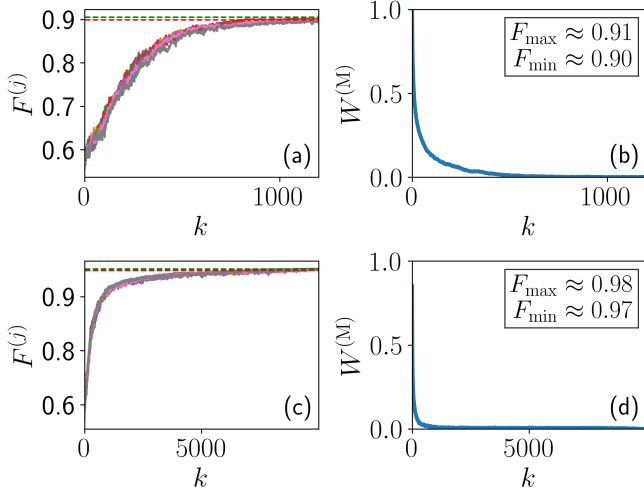


FIG. 3. Same as Fig. 2, but for the case of three qubits ( $d = 8$ ). All parameters are identical to those used for the two-qubit case shown in Fig. 2, except for the increased Hilbert-space dimension.

ously.

When comparing the results presented here with those of Ref. [34] for the two-qubit case—the only one reported in that work—we observe that our algorithm achieves slightly higher fidelities across the full set of states. In the method of Ref. [34], some states are obtained with high accuracy, whereas others perform less favorably. This difference may be related to the fact that, in our approach, all states are treated on equal footing, while in the other method some states are obtained earlier in the sequential procedure. Another important aspect is that the number of iterations required for convergence is significantly smaller in our algorithm. In fact, when attempting to extend the protocol of Ref. [34] to three qubits, we found that convergence could be reached for some states, but obtaining convergence for all of them within a reasonable number of iterations proved challenging. In the next section, we apply the present algorithm to paradigmatic physical Hamiltonians, including the transverse-field Ising model (TFIM) and Richardson pairing Hamiltonian, to demonstrate its robustness and versatility beyond random systems.

### III. APPLICATION TO PHYSICAL SYSTEMS

In this section, we apply the method proposed above to two types of physical systems. The first, relevant to condensed matter physics, is the TFIM [74], which is commonly used to describe condensed matter systems and has also frequently served as a benchmark in quantum computing studies. Assuming only nearest-neighbor in-

teractions, the TFIM Hamiltonian can be written as [75]

$$H = -J \sum_{j=0}^{N-2} Z_j Z_{j+1} - h \sum_{j=0}^{N-1} X_j - K \sum_{j=0}^{N-2} X_j Z_{j+1}. \quad (18)$$

Here  $j \in \{0, \dots, N-1\}$  labels the  $N$  qubits,  $\{X_j, Y_j, Z_j\}$  denotes the set of associated Pauli matrices, and  $J$ ,  $h$ , and  $K$  are parameters with dimensions of energy.

The second model is the pairing Hamiltonian, also known as the Richardson Hamiltonian [76–79], which is often used to describe small superfluid systems, such as atomic nuclei. Here, we focus on the specific subspace with seniority zero (no pair breaking) and employ the pair-to-qubit encoding technique. In this encoding, the Hamiltonian reads [80]

$$H = \sum_{j=0}^{N-1} \left( \varepsilon_j - \frac{g}{2} \right) (I - Z_j) - \frac{g}{2} \sum_{j=1}^{N-1} \sum_{k=0}^{j-1} (X_j X_k + Y_j Y_k), \quad (19)$$

where  $\varepsilon_i$  and  $g$  are parameters with dimensions of energy, representing the single-particle energy levels and the interaction strength, respectively. In the following, we assume equidistant single-particle levels,  $\varepsilon_j = j\Delta\varepsilon$ , with  $\Delta\varepsilon$  denoting the energy spacing between levels.

Both models share the feature of exhibiting a quantum phase transition (QPT). In the TFIM, the transition occurs from ordered to disordered spin orientations, whereas in the pairing model it corresponds to a  $U(1)$  symmetry-breaking transition associated with particle-number conservation. Beyond the presence of QPTs, however, the two models represent very different physical systems, both in nature and in their quantum information content. Specifically, the TFIM involves only local interactions and is generally characterized by area-law entanglement, while the pairing model features all-to-all interactions and is typically associated with volume-law entanglement.

In the following subsections, we present the results obtained by implementing the algorithm on these two models. As in the case of the random Hamiltonians analyzed earlier, both the TFIM and pairing Hamiltonians will first be rescaled to dimensionless form according to Eq. (17). In addition, at each iteration the dimensionless time  $\tilde{\tau}$  will be randomly sampled within an appropriate interval, as described in Sec. IID, to prevent convergence to spurious eigenstates.

#### A. Application to the TFIM model

To evaluate the efficiency of the algorithm when applied to the TFIM Hamiltonian, we first analyze the mean fidelities  $F_k^{(j)}$  shown in the left panels of Fig. 4. In these simulations, following the procedure described in Sec. IID, the dimensionless propagation time  $\tilde{\tau}$  is randomly sampled at each iteration within the interval  $[0, 600]$  to prevent convergence to spurious stationary states. In addition, the reset mechanism outlined



in Sec. II E is applied with parameters  $w_{\text{th}} = 0.005$  and  $w_r = 0.01$ . As shown in the figure, a slight decrease in both the maximum and minimum fidelities is observed as the number of qubits increases:  $F_{\text{max}} = 0.996$  and  $F_{\text{min}} = 0.993$  for two qubits,  $F_{\text{max}} = 0.988$  and  $F_{\text{min}} = 0.985$  for three qubits, and  $F_{\text{max}} = 0.968$  and  $F_{\text{min}} = 0.961$  for four qubits. This reduction in fidelities is expected, as increasing the number of qubits exponentially enlarges the dimension of the Hilbert space, making the exploration of the state space more challenging.

In the right panels of Fig. 4, we show, for each stochastic realization and for each computational basis state  $|j\rangle$ , the expectation values of the dimensionless TFIM Hamiltonian after the algorithm has converged, plotted against the index assigned to each realization. Specifically, for sufficiently large values of  $k$  such that the algorithm has effectively converged, we represent

$$\langle \tilde{H} \rangle_j = \langle j | D_k^\dagger \tilde{H} D_k | j \rangle \quad (20)$$

for each  $j \in \{0, \dots, d-1\}$  and for each realization, with  $\tilde{H}$  given by Eq. (17) using the TFIM Hamiltonian of Eq. (18). The solid horizontal lines indicate the exact values obtained via numerical diagonalization of the Hamiltonian. As can be seen, for the two- and three-qubit cases, the energies obtained in most realizations match the exact numerical values quite closely. This indicates that, in these cases, even a single realization is sufficient to approximate the full set of eigenstates. For four qubits, however, deviations become more noticeable, even though the fidelities are only slightly lower than in the smaller systems. This highlights that the computation of these energies is highly sensitive to small variations in the stationary states approximated by the algorithm.

In an effort to improve the estimation of the eigenenergies, we slightly modified the algorithm parameters. Specifically, in Fig. 5 we show the same results as in Fig. 4, but using a slightly higher reward parameter,  $r = 0.93$ . Compared to Fig. 4, this adjustment of the reward parameter generally leads to a modest improvement in both the fidelities and the estimated energies, indicating that small changes in the algorithm settings can enhance the accuracy of the results across all realizations.

It should be noted that the example discussed above represents an ideal benchmark, as the performance of the algorithm can be directly assessed by comparing its output with the exact stationary states and eigenenergies obtained through numerical diagonalization of the Hamiltonian. In more realistic scenarios, however, such exact information may not be available, making it impossible to evaluate the accuracy of the algorithm in a direct manner. Nevertheless, for the eigenenergies, it is still possible to identify, among all stochastic realizations of the algorithm, those providing the best approximations to the ground and highest-energy states. Specifically, the realizations yielding the smallest and largest expectation values of the Hamiltonian can be associated with these

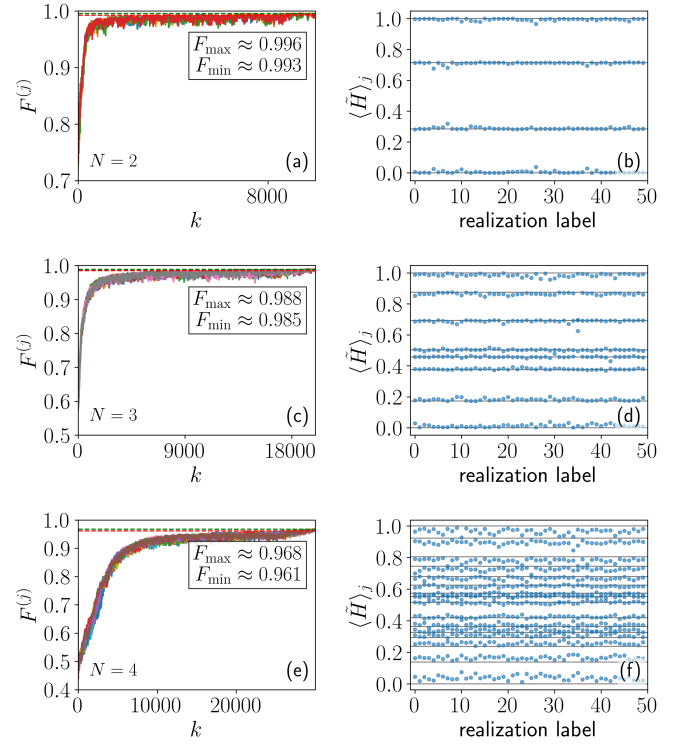


FIG. 4. Results of the algorithm applied to the TFIM Hamiltonian in Eq. (18) with  $J/h = 1$  and  $K/h = 0.5$ . The remaining parameters are  $r = 0.9$ ,  $p = 2/r$ ,  $w_{\text{th}} = 0.005$ ,  $w_r = 0.05$ , and  $N_r = 50$ . In each iteration, the dimensionless evolution time  $\tilde{\tau}$  is uniformly sampled from  $[0, 600]$ . The upper panels [(a) and (b)] correspond to the two-qubit case ( $d = 4$ ), the middle panels [(c) and (d)] to three qubits ( $d = 8$ ), and the lower panels [(e) and (f)] to four qubits ( $d = 16$ ). In the left panels [(a), (c), and (e)], the mean fidelities  $F_k^{(j)}$  associated with each computational-basis state  $|j\rangle$  are shown as a function of the iteration number  $k$ , with different colors indicating different states. Since many curves overlap, color differences may be difficult to distinguish. Dashed horizontal lines indicate the maximum and minimum fidelities,  $F_{\text{max}}$  and  $F_{\text{min}}$ . The right panels [(b), (d), and (f)] display, for each computational basis state  $|j\rangle$ , the dimensionless expectation values of the energies obtained after each stochastic realization has converged, plotted against the corresponding realization index. In these panels, horizontal solid lines represent the exact eigenenergies computed by numerical diagonalization.

two states, respectively. This identification is justified by the variational principle, which states that the expectation value of the Hamiltonian in any normalized state constitutes an upper bound to the exact ground-state energy. Likewise, applying the same argument to the inverted Hamiltonian,  $-H$ , the expectation value provides a lower bound to the highest-energy eigenvalue. Consequently, the realizations corresponding to the smallest and largest energy expectations offer the most accurate approximations to the ground and highest-energy states, respectively. Later, we will introduce a more general criterion based on evaluating the energy dispersions of the



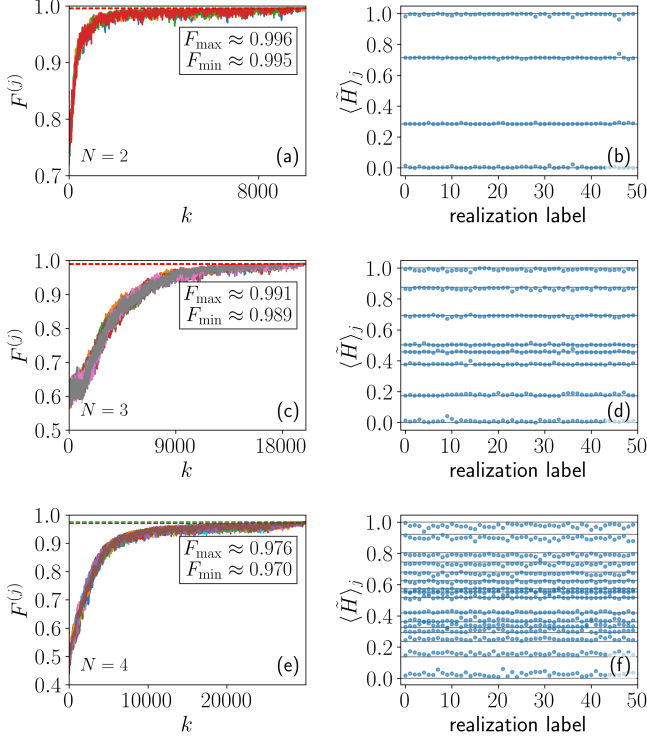


FIG. 5. Results analogous to those shown in Fig. 4, obtained for  $r = 0.93$ . All other parameters are identical to those used in Fig. 4.

states obtained from the different realizations of the algorithm.

### B. Application to the pairing model

We now consider the pairing Hamiltonian given by Eq. (19). This Hamiltonian describes a set of interacting particle pairs and is relevant for modeling the transition from normal to superfluid phases in small superconducting systems. Notably, while in the TFIM model the exact solution required the diagonalization of a  $d \times d$  matrix, here the Hamiltonian exhibits a block-diagonal structure in the computational qubit basis due to particle number conservation. Specifically, after applying the particle-to-qubit transformation leading to Eq. (19), particle number conservation is mapped onto the conservation of the Hamming weight. Starting from a given basis state  $|j\rangle$ , we can equivalently write it in the qubit basis as  $|[j]\rangle$ , where  $[j]$  denotes the binary representation of  $j$ . The Hamming weight is defined as the number of “1”s in  $[j]$ . In the pairing Hamiltonian, two basis states  $|j\rangle$  and  $|k\rangle$  have non-zero Hamiltonian matrix elements only if they share the same Hamming weight. Consequently, the block corresponding to Hamming weight  $N_{\text{Ham}}$  has dimension  $C_N^{N_{\text{Ham}}}$ , the binomial coefficient (“ $N$  choose  $N_{\text{Ham}}$ ”), which counts the number of ways to choose  $N_{\text{Ham}}$  objects out of  $N$ . This block structure

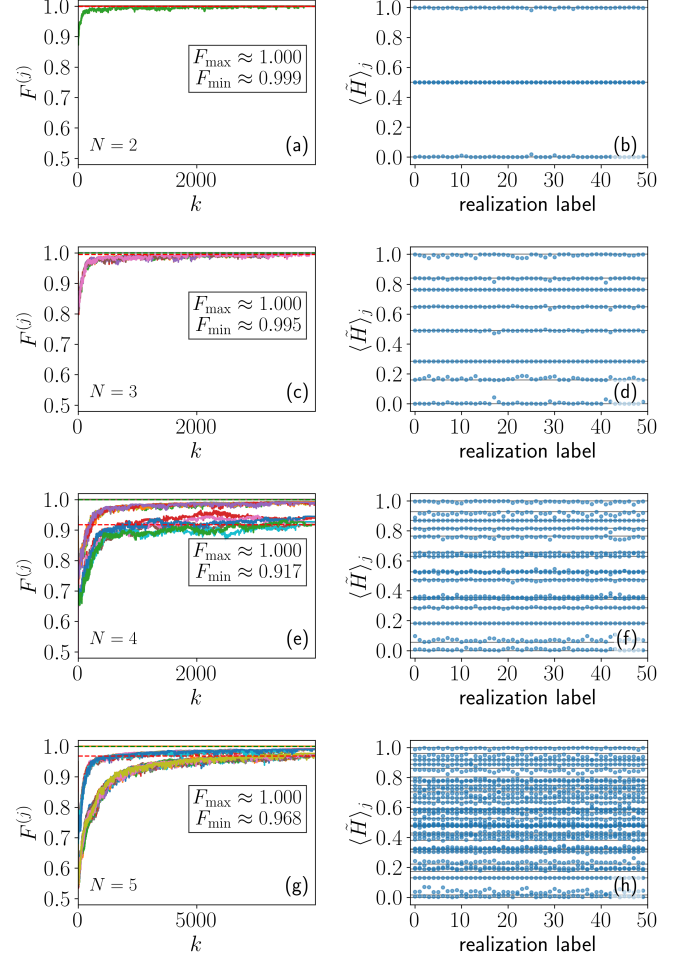


FIG. 6. Results of the algorithm applied to the pairing Hamiltonian in Eq. (19) with  $\varepsilon_j = j\Delta\varepsilon$  and  $g/\Delta\varepsilon = 1$ . The remaining parameters are  $r = 0.9$ ,  $p = 2/r$ ,  $w_{\text{th}} = 0.005$ ,  $w_r = 0.05$ , and  $N_r = 50$ . In each iteration, the dimensionless evolution time  $\tilde{\tau}$  is uniformly sampled from  $[0, 600]$ . Panels (a)–(b) show the two-qubit case ( $d = 4$ ), (c)–(d) the three-qubit case ( $d = 8$ ), (e)–(f) the four-qubit case ( $d = 16$ ), and (g)–(h) the five-qubit case ( $d = 32$ ). The left panels show the mean fidelities  $F_k^{(j)}$  associated with each computational-basis state  $|j\rangle$  as a function of the iteration number  $k$ , with different colors indicating different states. Since many curves overlap, color differences may be difficult to distinguish. Dashed horizontal lines indicate the maximum and minimum fidelities,  $F_{\text{max}}$  and  $F_{\text{min}}$ . The right panels display, for each computational-basis state  $|j\rangle$ , the dimensionless expectation values of the energies obtained after each stochastic realization has converged, plotted against the corresponding realization index. In these panels, horizontal solid lines represent the exact eigenenergies computed by numerical diagonalization.

is commonly exploited in classical computations to access eigenstates by diagonalizing the Hamiltonian within each block. In the context of machine learning, however, exploiting problem symmetries can introduce additional challenges. Here, one of our objectives is to test whether the RL method can converge to the correct energies with-

out explicitly enforcing these symmetries during learning.

Figure 6 presents results analogous to those in Fig. 4, but for the pairing Hamiltonian instead of the TFIM model. The left panels display the mean fidelities for each computational-basis state  $|j\rangle$ , averaged over all stochastic realizations, as a function of the iteration number. The right panels show, for each realization and computational-basis state, the expectation values of the dimensionless pairing Hamiltonian obtained after the algorithm has converged, compared with the exact eigenenergies from numerical diagonalization. The results are shown for increasing numbers of qubits, from  $N = 2$  to  $N = 5$ . The coupling strength parameter is set to  $g/\Delta\varepsilon = 1$ , corresponding to the strong-coupling regime that gives rise to the superfluid phase. Accordingly, the eigenstates generally involve complex superpositions of the original basis states. However, there are two particular computational basis states,  $|0\rangle$  and  $|d-1\rangle$ , which correspond to Hamming weights  $N_{\text{Ham}} = 0$  and  $N_{\text{Ham}} = N$ , respectively, and are exact eigenstates of the pairing Hamiltonian. Consequently, two horizontal lines with fidelity equal to 1 appear in the left panels of the figure, and the maximum fidelity is therefore always 1 in those plots.

Although the properties of the pairing Hamiltonian differ substantially from those of the TFIM model, we again observe that, after a reasonable number of iterations, the algorithm achieves very good fidelities in most cases, with the exception of the four-qubit system. As seen in the left panels, the mean fidelities remain high overall, although the minimum fidelity tends to decrease as the number of qubits increases. This trend is expected, since the dimension of the Hilbert space grows exponentially with the system size, making the exploration of the state space increasingly demanding. This behavior is consistent with the results shown in the right panels, where the expectation values of the dimensionless Hamiltonian are generally very close to the exact eigenenergies across nearly all realizations. An exception occurs for  $N = 4$ , where two distinct groups of fidelities appear in panel (e): one set converging close to unity and another slightly above 0.9, resulting in a minimum fidelity smaller than that for  $N = 5$ . This apparent inconsistency with the accurate reproduction of the eigenenergies [panel (f)] can be explained by the presence of an exact degeneracy in one of the eigenenergies. Within the corresponding degenerate subspace, the eigenstates are not uniquely defined, and the algorithm may converge to different linear combinations than those obtained from direct diagonalization, leading to slightly reduced fidelities. Nevertheless, the method performs well for  $N = 4$ , successfully recovering the correct set of eigenenergies.

For  $N = 5$  [panels (g) and (h)], three distinct groups of curves appear, all reaching high fidelities but requiring different numbers of iterations to converge. In the final iterations, the expectation values of the dimensionless Hamiltonian remain reasonably close to the exact energies, though the comparison becomes less straightforward

due to the increased density of states.

The emergence of multiple convergence rates arises from the symmetry structure of the pairing Hamiltonian, which organizes the Hilbert space into blocks of different sizes according to particle-number conservation. In the qubit representation, these blocks correspond to subspaces with fixed Hamming weight. For five qubits, there are two blocks of size 1 ( $N_{\text{Ham}} = 0$  and  $N_{\text{Ham}} = 5$ ), two of size 5 ( $N_{\text{Ham}} = 1$  and  $N_{\text{Ham}} = 4$ ), and two of size 10 ( $N_{\text{Ham}} = 2$  and  $N_{\text{Ham}} = 3$ ). The cases  $N_{\text{Ham}} = 0$  and  $N_{\text{Ham}} = 5$  correspond to the exact eigenstates  $|0\rangle$  and  $|31\rangle$ , explaining the horizontal lines with fidelity one in panel (g). The groups of curves that converge faster correspond to the blocks of size 5, while the slower convergence is associated with the larger blocks of size 10. This behavior reflects that the number of iterations required for convergence increases with the dimension of the explored subspace. It is worth noting that the algorithm manages to converge within all symmetry sectors without being provided with any prior information about the underlying symmetries.

### C. Symmetry-restricted RL

We showed in the previous section that the RL protocol introduced here is able to capture the symmetries of a Hamiltonian without prior knowledge of them. Alternatively, if the symmetries of the system are known in advance, the RL protocol can be directly applied within a given symmetry sector. Exploiting such symmetries can effectively extend the applicability of the RL technique to systems with a larger number of qubits, since each sector has a lower dimension than the full Hilbert space.

To illustrate the above-mentioned procedure, we first consider the case  $N = 5$  shown in panels (g) and (h) of Fig. 6. To account for particle-number conservation, instead of performing the RL in the full Hilbert space of dimension  $d = 2^5 = 32$ , we apply the RL algorithm separately within sectors of states  $|j\rangle$  whose Hamming weight takes a fixed value between 0 and 5. Panels (a) and (b) of Fig. 7 show the results obtained from the independent symmetry-restricted calculations corresponding to each Hamming-weight sector:  $N_{\text{Ham}} = 0$  ( $d = C_5^0 = 1$ ),  $N_{\text{Ham}} = 1$  ( $d = C_5^1 = 5$ ),  $N_{\text{Ham}} = 2$  ( $d = C_5^2 = 10$ ),  $N_{\text{Ham}} = 3$  ( $d = C_5^3 = 10$ ),  $N_{\text{Ham}} = 4$  ( $d = C_5^4 = 5$ ), and  $N_{\text{Ham}} = 5$  ( $d = C_5^5 = 1$ ). Since the dimension of each of these sectors is smaller than that of the full Hilbert space, the application of the algorithm becomes less demanding. This figure also confirms our previous interpretation that the different convergence scales observed earlier originate from the different dimensions of the symmetry sectors.

In order to highlight the advantages of this approach, panels (c) and (d) of Fig. 7 show the results of applying the symmetry-restricted RL algorithm to a larger system with  $N = 6$  qubits. The fact that some sectors display lower fidelity values does not necessarily indicate a failure of the algorithm. This behavior is analogous to that

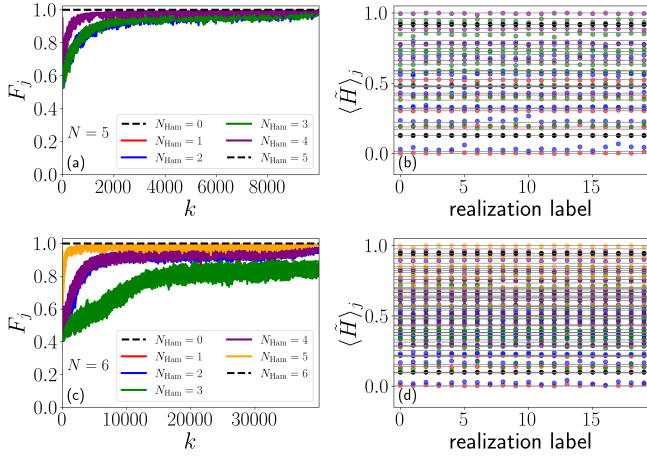


FIG. 7. Results of the symmetry-restricted RL algorithm applied to the pairing Hamiltonian. The left panels display, for each computational-basis state, the mean fidelities as a function of the iteration number, while the right panels show, for each state and realization, the dimensionless expectation values of the energies obtained after the stochastic process has converged. Different colors represent the results obtained for the different Hamming-weight sectors. Panels (a) and (b) correspond to the five-qubit system ( $N = 5$ ) shown in panels (g) and (h) of Fig. 6, but with the RL applied independently within each Hamming-weight sector. Panels (c) and (d) show analogous results for a six-qubit system ( $N = 6$ ). All parameter values are the same as in Fig. 6, except for the number of realizations, which is  $N_r = 20$  in this case.

observed for the four-qubit case in panel (e) of Fig. 6. In the present case, several eigenenergies are degenerate, and within each degenerate subspace, the eigenstates are not uniquely defined. As a result, the algorithm may converge to linear combinations of these degenerate eigenstates that differ from those obtained by direct diagonalization, leading to reduced fidelities without implying incorrect convergence.

#### D. Post-selection of “good” final states

An interesting aspect illustrated by the previous figures is the dispersion of the final expectation values of the dimensionless Hamiltonian observed among different realizations. Even though the mean fidelities are high, these expectation values may differ noticeably from the exact eigenenergies from one realization to another. Similarly, the fidelities of the final states also fluctuate between realizations, despite the relatively high mean value. A natural question that arises is how to determine whether the results obtained in a specific realization are good or bad. The fidelities associated with that realization cannot serve as a criterion, since they could only be computed if the exact eigenstates were already known. Moreover, as previously discussed, low fidelities do not necessarily imply poor convergence to the eigenstates in the

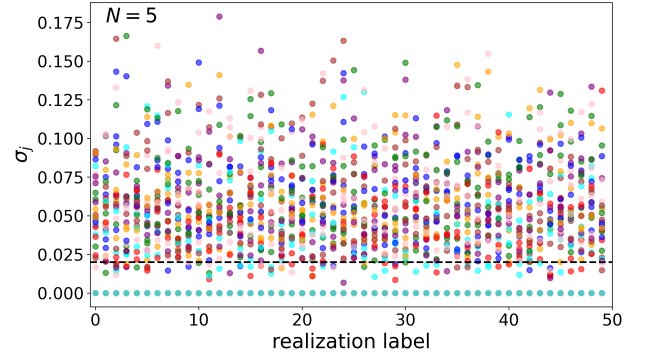


FIG. 8. Illustration of the evaluation of the energy fluctuations defined in Eq. (21). This quantity is computed for each final state and plotted as a function of the realization number. The results correspond to the final states obtained for the pairing Hamiltonian with  $N = 5$  qubits, as in panels (g) and (h) of Fig. 6 and panels (a) and (b) of Fig. 7. All parameter values are the same as in those figures. The horizontal dashed line indicates the threshold value  $\sigma_{th} = 0.02$ , which is used to select the final states shown in Fig. 10.

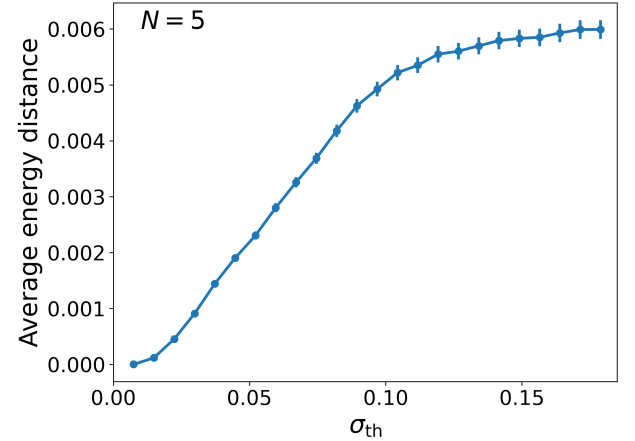


FIG. 9. Average distance between the expectation value  $\langle \tilde{H} \rangle_j$  and the closest eigenenergy of the dimensionless pairing Hamiltonian, shown as a function of the threshold value  $\sigma_{th}$ . The average is computed over all final states and stochastic realizations in Fig. 8 that satisfy the condition  $\sigma_j \leq \sigma_{th}$ . Error bars indicate statistical deviations across the selected realizations and states.

presence of degeneracies. As an alternative, one can consider the energy fluctuations, defined for a given state and realization as

$$\sigma_j = \sqrt{\langle j | D_k^\dagger \tilde{H}^2 D_k | j \rangle - \langle j | D_k^\dagger \tilde{H} D_k | j \rangle^2}, \quad (21)$$

where  $k$  is chosen large enough for the algorithm to have converged. Indeed, this quantity is always non-negative and vanishes only when the state is an exact eigenstate.

Figure 8 provides a quantitative illustration of the final

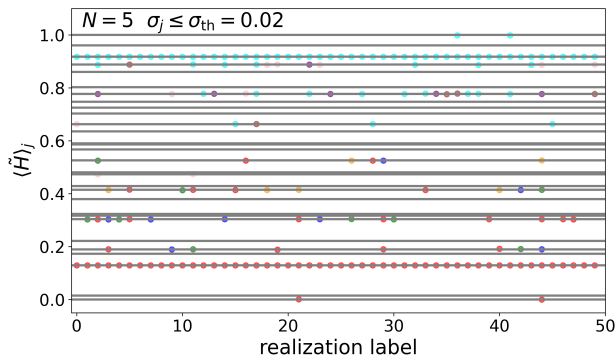


FIG. 10. Effect of introducing a fluctuation threshold on the final estimated energies. Only states satisfying  $\sigma \leq \sigma_{th} = 0.02$  are shown in the energy-versus-realization plot. The results correspond to the  $N = 5$  pairing Hamiltonian, using the same data as in Fig. 8. The effect of this post-selection can be appreciated by comparison with panel (f) of Fig. 6, where all final states are included.

energy fluctuations computed for the states obtained at the end of the learning process with the pairing Hamiltonian for the case of  $N = 5$  qubits. The results correspond to the same parameter set used in panels (g) and (h) of Fig. 6 and panels (a) and (b) of Fig. 7. The figure clearly shows that the fluctuation values exhibit a rather wide dispersion among the different computational-basis states, even within a single realization.

It is to be expected that, for a given state, smaller energy fluctuations indicate a closer convergence to an eigenstate of the Hamiltonian. To further illustrate this, Fig. 9 shows the average distance between the expectation value  $\langle \tilde{H} \rangle_j$  and the closest eigenenergy of the dimensionless pairing Hamiltonian as a function of a threshold value  $\sigma_{th}$ . The average is taken over all final states and realizations from Fig. 8 that satisfy  $\sigma_j \leq \sigma_{th}$ , while error bars denote statistical deviations across the selected realizations and states. As observed, the smaller the threshold  $\sigma_{th}$ , the smaller the distance between the expectation values and the exact eigenenergies. This observation motivates the introduction of a fluctuation threshold  $\sigma_{th}$  to post-select the most reliable final states.

To illustrate that energy fluctuations can be used as a post-selection criterion to identify “good” final states, we present in Fig. 10 the final expectation values of the dimensionless Hamiltonian, retaining only those states whose energy fluctuations fall below a chosen threshold,  $\sigma_{th} = 0.02$ . Even when the RL algorithm has not fully converged for all states within a given realization, this selection allows us to isolate the most reliable results. The expectation values obtained from the selected final states show a much better agreement with the exact eigenenergies than those displayed in panel (f) of Fig. 6, where all final states were included. This confirms that imposing a fluctuation-based selection efficiently filters out poorly converged states while preserving those that faithfully

reproduce the true eigenstates of the Hamiltonian.

#### IV. CONCLUSIONS

In this work, we have introduced a RL-based algorithm designed to identify the unitary transformation that maps the computational basis onto the set of pure states that remain invariant under the action of a given quantum operation. A key feature of this approach is its ability to determine all these states simultaneously and the way it constructs the global unitary transformation iteratively from two-qudit operations guided by a reinforcement strategy.

To benchmark the algorithm, we considered the problem of finding the eigenstates of specific Hamiltonians. For pseudo-random Hamiltonians involving two- and three-qubit systems, the method achieved high fidelities across multiple stochastic realizations, confirming its robustness and accuracy. We then applied the algorithm to two physical models: the TFIM and the pairing Hamiltonian. In the TFIM case, high fidelities were maintained as the number of qubits increased from two to four, demonstrating the scalability of the approach. For the pairing Hamiltonian, reliable convergence was obtained for systems of up to six qubits, although some states exhibited slightly lower fidelities due to the larger Hilbert-space dimension and the presence of degeneracies.

A remarkable feature of the algorithm is its ability to uncover the underlying symmetries of the system without prior knowledge. This manifests in the appearance of different convergence time scales associated with symmetry sectors of varying sizes. Moreover, by explicitly restricting the learning process to subspaces of fixed Hamming weight, we verified that the algorithm correctly identifies the eigenstates within each symmetry block. Such restriction also makes the computation significantly less demanding, as the effective dimensionality of the learning space is reduced.

We also introduced a post-selection criterion based on the energy fluctuations of the final states, which allows one to identify well-converged (“good”) results even when the RL process has not fully converged for all states. Applying a threshold on these fluctuations efficiently filters out poorly converged states, yielding a set of results that almost perfectly reproduce the exact eigenenergies.

Although in this work we have focused on unitary quantum operations, the structure of the algorithm is not limited to this case. Its formulation naturally extends to non-unitary dynamics, making it a promising tool for investigating dissipative or open quantum systems. Future work will explore this extension and, in particular, the case where the fixed points of the quantum operation are not pure states forming an orthonormal basis.

## ACKNOWLEDGMENTS

The authors acknowledge project PID2022-136228NB-C22, funded by MCIN/AEI/10.13039/501100011033 and by “ERDF A way of making Europe”, EU, which has specifically covered the publication costs of this article. Furthermore, this work has been partially financially supported by the Ministry of Economic Affairs and Digital Transformation of the Spanish Government through the QUANTUM ENIA project call—Quantum Spain project, and by the European Union through the Recovery, Transformation and Resilience Plan—NextGenerationEU within the framework of the “Digital Spain 2026 Agenda”. It has also been co-financed by the EU, Ministerio de Hacienda y Función

Pública, FEDER, and Junta de Andalucía (project SOL2024-31833). This work has also received financial support from the CNRS through the AIQI-IN2P3 project. Support from the CNRS-IRN ASTRANUCAP and the Erasmus program project 2024-1-ES01-KA131-HED-000219547 for the visit of M.L.O.-A. and J.C.-P. to D.L. is gratefully acknowledged. This work is part of the HQI initiative ([www.hqi.fr](http://www.hqi.fr)) and is supported by France 2030 under the French National Research Agency award number “ANR-22-PNQC-0002”.

## DATA AVAILABILITY

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

- 
- [1] M. A. Nielsen and I. L. Chuang, *Quantum Computing and Quantum Information* (Cambridge University Press, Cambridge, 2000).
  - [2] M. Schuld, I. Sinayskiy, and F. Petruccione, An introduction to quantum machine learning, *Contemp. Phys.* **56**, 172 (2015).
  - [3] G. Torlai, G. Mazzola, J. Carrasquilla, M. Troyer, R. Melko, and G. Carleo, Neural-network quantum state tomography, *Nat. Phys.* **14**, 447 (2018).
  - [4] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, Quantum machine learning, *Nature* **549**, 195 (2017).
  - [5] G. Liu, M. Chen, Y.-X. Liu, D. Layden, and P. Cappelaro, Repetitive readout enhanced by machine learning, *Mach. Learn.: Sci. Technol.* **1**, 015003 (2020).
  - [6] R. Heese, T. Gerlach, S. Mücke, S. Müller, M. Jakobs, and N. Piatkowski, Explaining quantum circuits with shapley values: towards explainable quantum machine learning, *Quantum Mach. Intell.* **7**, 10.1007/s42484-025-00254-8 (2025).
  - [7] T. Haug, C. N. Self, and M. S. Kim, Quantum machine learning of large datasets using randomized measurements, *Mach. Learn.: Sci. Technol.* **4**, 015005 (2023).
  - [8] E. Perrier, D. Tao, and C. Ferrie, Quantum geometric machine learning for quantum circuits and control, *New J. Phys.* **22**, 103056 (2020).
  - [9] F. A. Cárdenas-López, M. Sanz, J. C. Retamal, and E. Solano, Enhanced quantum synchronization via quantum machine learning, *Adv. Quantum Technol.* **2**, 1800076 (2019).
  - [10] S. Saini, P. Khosla, M. Kaur, and G. Singh, Quantum driven machine learning, *Int. J. Theor. Phys.* **59**, 4013 (2020).
  - [11] N. Innan, O. I. Siddiqui, S. Arora, T. Ghosh, Y. P. Koçak, D. Paragas, A. A. O. Galib, M. A.-Z. Khan, and M. Ben-nai, Quantum state tomography using quantum machine learning, *Quantum Mach. Intell.* **6** (2024).
  - [12] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction.*, 2nd ed., Adaptive Computation and Machine Learning (MIT Press, Cambridge, 2018).
  - [13] S. Shrapnel, F. Costa, and G. Milburn, Quantum Markovianity as a supervised learning task, *Int. J. Quantum Info.* **16**, 1840010 (2018).
  - [14] I. Agresti, N. Viggianiello, F. Flamini, N. Spagnolo, A. Crespi, R. Osellame, N. Wiebe, and F. Sciarrino, Pattern recognition techniques for boson sampling validation, *Phys. Rev. X* **9**, 011013 (2019).
  - [15] A. Youssef, R. J. Chapman, A. Peruzzo, C. Ferrie, and M. Tomamichel, Modeling and control of a reconfigurable photonic circuit using deep learning, *Quantum Sci. Technol.* **5**, 025001 (2020).
  - [16] I. A. Luchnikov, S. V. Vintskevich, D. A. Grigoriev, and S. N. Filippov, Machine Learning Non-Markovian Quantum Dynamics, *Phys. Rev. Lett.* **124**, 140502 (2020).
  - [17] A. A. Melnikov, L. E. Fedichkin, R.-K. Lee, and A. Alodjants, Machine learning transfer efficiencies for noisy quantum walks, *Adv. Quantum Technol.* **3**, 1900115 (2020).
  - [18] B. Lucian, B. Robert, S. B. De, and E. Damien, *Reinforcement Learning and Dynamic Programming Using Function Approximators*, 1st ed. (CRC Press Inc., Boca Raton, FL, USA, 2010).
  - [19] F. Agostinelli, G. Hocquet, S. Singh, and P. Baldi, From reinforcement learning to deep reinforcement learning: An overview, in *Braverman Readings in Machine Learning. Key Ideas from Inception to Current State: International Conference Commemorating the 40th Anniversary of Emmanuil Braverman's Decease, Boston, MA, USA, April 28-30, 2017, Invited Talks*, edited by L. Rozonoer, B. Mirkin, and I. Muchnik (Springer International Publishing, Cham, 2018) pp. 298–328.
  - [20] D. Dong, C. Chen, H. Li, and T.-J. Tarn, Quantum reinforcement learning, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **38**, 1207 (2008).
  - [21] G. D. Paparo, V. Dunjko, A. Makmal, M. A. Martin-Delgado, and H. J. Briegel, Quantum speedup for active learning agents, *Phys. Rev. X* **4**, 031002 (2014).
  - [22] V. Dunjko, J. M. Taylor, and H. J. Briegel, Quantum-enhanced machine learning, *Phys. Rev. Lett.* **117**, 130501 (2016).



- [23] 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).
- [24] M. Bukov, Reinforcement learning for autonomous preparation of floquet-engineered states: Inverting the quantum kapitza oscillator, *Phys. Rev. B* **98**, 224305 (2018).
- [25] T. Fösel, P. Tighineanu, T. Weiss, and F. Marquardt, Reinforcement learning with neural networks for quantum feedback, *Phys. Rev. X* **8**, 031084 (2018).
- [26] J. Kaldari, S. Tariq, S. Al-Kuwari, S. Y.-C. Chen, S. Chatzinotas, and H. Shin, Quantum reinforcement learning: Recent advances and future directions (2025), arXiv:2510.14595 [quant-ph].
- [27] C. Jiang, Y. Pan, Z.-G. Wu, Q. Gao, and D. Dong, Robust optimization for quantum reinforcement learning control using partial observations, *Phys. Rev. A* **105**, 062443 (2022).
- [28] E. Cherrat, I. Kerenidis, and A. Prakash, Quantum reinforcement learning via policy iteration, *Quantum Mach. Intell.* **5**, DOI10.1007/s42484-023-00116-1 (2023).
- [29] Y. Liu, C. Xu, and S. Jin, Reinforcement learning for continuous control: A quantum normalized advantage function approach, in *2023 IEEE International Conference on Quantum Software (QSW)* (2023) pp. 83–87.
- [30] G. Franceschetto and A. Ricou, Demonstration of quantum projective simulation on a single-photon-based quantum computer, *Phys. Rev. A* **110**, 062613 (2024).
- [31] M. L. Olivera-Atencio, L. Lamata, M. Morillo, and J. Casado-Pascual, Quantum reinforcement learning in the presence of thermal dissipation, *Phys. Rev. E* **108**, 014128 (2023).
- [32] M. L. Olivera-Atencio, L. Lamata, and J. Casado-Pascual, Benefits of open quantum systems for quantum machine learning, *Adv. Quantum Technol.*, 2300247 (2023).
- [33] M. L. Olivera-Atencio, L. Lamata, and J. Casado-Pascual, Impact of amplitude and phase damping noise on quantum reinforcement learning: challenges and opportunities, *Eur. Phys. J. Spec. Top.* (2025).
- [34] F. Albarrán-Arriagada, J. C. Retamal, E. Solano, and L. Lamata, Reinforcement learning for semi-autonomous approximate quantum eigensolver, *Mach. Learn.: Sci. Technol.* **1**, 015002 (2020).
- [35] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, The theory of variational hybrid quantum-classical algorithms, *New J. Phys.* **18**, 023023 (2016).
- [36] S. McArdle, S. Endo, A. Aspuru-Guzik, S. C. Benjamin, and X. Yuan, Quantum computational chemistry, *Rev. Mod. Phys.* **92**, 015003 (2020).
- [37] M. Motta and J. E. Rice, Emerging quantum computing algorithms for quantum chemistry, *WIREs Comput. Mol. Sci.* **12**, e1580 (2022).
- [38] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, and P. J. Coles, Variational quantum algorithms, *Nat. Rev. Phys.* **3**, 625 (2021).
- [39] J. Tilly, H. Chen, S. Cao, D. Picozzi, K. Setia, Y. Li, E. Grant, L. Wossnig, I. Rungger, G. H. Booth, and J. Tennyson, The variational quantum eigensolver: A review of methods and best practices, *Phys. Rep.* **986**, 1 (2022).
- [40] M. Larocca, S. Thanasilp, S. Wang, K. Sharma, J. Biamonte, P. J. Coles, L. Cincio, J. R. McClean, Z. Holmes, and M. Cerezo, Barren plateaus in variational quantum computing, *Nat. Rev. Phys.* **7**, 174 (2025).
- [41] J. R. McClean, M. E. Kimchi-Schwartz, J. Carter, and W. A. de Jong, Hybrid quantum-classical hierarchy for mitigation of decoherence and determination of excited states, *Phys. Rev. A* **95**, 042308 (2017).
- [42] H. R. Grimsley, S. E. Economou, E. Barnes, and N. J. Mayhall, An adaptive variational algorithm for exact molecular simulations on a quantum computer, *Nat. Commun.* **10**, 3007 (2019).
- [43] K. M. Nakanishi, K. Mitarai, and K. Fujii, Subspace-search variational quantum eigensolver for excited states, *Phys. Rev. Research* **1**, 033062 (2019).
- [44] R. M. Parrish, E. G. Hohenstein, P. L. McMahon, and T. J. Martínez, Quantum computation of electronic transitions using a variational quantum eigensolver, *Phys. Rev. Lett.* **122**, 230401 (2019).
- [45] O. Higgott, D. Wang, and S. Brierley, Variational quantum computation of excited states, *Quantum* **3**, 156 (2019).
- [46] T. Jones, S. Endo, S. McArdle, X. Yuan, and S. C. Benjamin, Variational quantum algorithms for discovering Hamiltonian spectra, *Phys. Rev. A* **99**, 062304 (2019).
- [47] M. Motta, C. Sun, A. T. K. Tan, M. J. O’Rourke, E. Ye, A. J. Minnich, F. G. S. L. Brandão, and G. K.-L. Chan, Determining eigenstates and thermal states on a quantum computer using quantum imaginary time evolution, *Nat. Phys.* **16**, 205 (2020).
- [48] W. J. Huggins, J. Lee, U. Baek, B. O’Gorman, and K. B. Whaley, A non-orthogonal variational quantum eigensolver, *New J. Phys.* **22**, 073009 (2020).
- [49] N. H. Stair, R. Huang, and F. A. Evangelista, A Multireference Quantum Krylov Algorithm for Strongly Correlated Electrons, *J. Chem. Theory Comput.* **16**, 2236 (2020).
- [50] K. Seki and S. Yunoki, Quantum power method by a superposition of time-evolved states, *PRX Quantum* **2**, 010333 (2021).
- [51] Y. S. Yordanov, V. Armaos, C. H. W. Barnes, and D. R. M. Arvidsson-Shukur, Qubit-excitation-based adaptive variational quantum eigensolver, *Commun. Phys.* **4**, 1 (2021).
- [52] H. L. Tang, V. O. Shkolnikov, G. S. Barron, H. R. Grimsley, N. J. Mayhall, E. Barnes, and S. E. Economou, Qubit-adapt-vqe: An adaptive algorithm for constructing hardware-efficient ansätze on a quantum processor, *PRX Quantum* **2**, 020310 (2021).
- [53] F. Zhang, N. Gomes, Y. Yao, P. P. Orth, and T. Iadecola, Adaptive variational quantum eigensolvers for highly excited states, *Phys. Rev. B* **104**, 075159 (2021).
- [54] Y. S. Yordanov, C. H. W. Barnes, and D. R. M. Arvidsson-Shukur, Molecular-excited-state calculations with the qubit-excitation-based adaptive variational quantum eigensolver protocol, *Phys. Rev. A* **106**, 032434 (2022).
- [55] C. L. Cortes and S. K. Gray, Quantum Krylov subspace algorithms for ground- and excited-state energy estimation, *Phys. Rev. A* **105**, 022417 (2022).
- [56] E. N. Epperly, L. Lin, and Y. Nakatsukasa, A theory of quantum subspace diagonalization, *SIAM J. Matrix Anal. Appl.* **43**, 1263 (2022).



- [57] E. A. R. Guzman and D. Lacroix, Accessing ground-state and excited-state energies in a many-body system after symmetry restoration using quantum computers, *Phys. Rev. C* **105**, 024324 (2022).
- [58] M. Q. Hlatshwayo, Y. Zhang, H. Wibowo, R. LaRose, D. Lacroix, and E. Litvinova, Simulating excited states of the Lipkin model on a quantum computer, *Phys. Rev. C* **106**, 024319 (2022).
- [59] Y. Wang and D. A. Mazziotti, Electronic excited states from a variance-based contracted quantum eigensolver, *Phys. Rev. A* **108**, 022814 (2023).
- [60] Y. Shen, K. Klymko, J. Sud, D. B. Williams-Young, W. A. de Jong, and N. M. Tubman, Real-Time Krylov Theory for Quantum Computing Algorithms, *Quantum* **7**, 1066 (2023).
- [61] M. Zheng, B. Peng, N. Wiebe, A. Li, X. Yang, and K. Kowalski, Quantum algorithms for generator coordinate methods, *Phys. Rev. Res.* **5**, 023200 (2023).
- [62] S. E. Smart, D. M. Welakuh, and P. Narang, Many-body excited states with a contracted quantum eigensolver, *J. Chem. Theory Comput.* **20**, 3580 (2024).
- [63] M. Motta, W. Kirby, I. Liepuoniute, K. J. Sung, J. Cohn, A. Mezzacapo, K. Klymko, N. Nguyen, N. Yoshioka, and J. E. Rice, Subspace methods for electronic structure simulations on quantum computers, *Electron. Struct.* **6**, 013001 (2024).
- [64] M. R. Hirsbrunner, J. W. Mullinax, Y. Shen, D. B. Williams-Young, K. Klymko, R. V. Beeumen, and N. M. Tubman, Diagnosing local minima and accelerating convergence of variational quantum eigensolvers with quantum subspace techniques (2024), arXiv:2404.06534 [quant-ph].
- [65] M. Zheng, B. Peng, A. Li, X. Yang, and K. Kowalski, Unleashed from constrained optimization: Quantum computing for quantum chemistry employing generator coordinate inspired method, *NPJ Quantum Inf.* **10**, 1 (2024).
- [66] Y. O. Nakagawa, M. Kamoshita, W. Mizukami, S. Sudo, and Y. ya Ohnishi, Adapt-qsci: Adaptive construction of an input state for quantum-selected configuration interaction, *J. Chem. Theory Comput.* **20**, 10817 (2024).
- [67] L. Ding, C.-L. Hong, and C. Schilling, Ground and excited states from ensemble variational principles, *Quantum* **8**, 1525 (2024).
- [68] Y. Beaujeault-Taudière and D. Lacroix, Solving the Lipkin model using quantum computers with two qubits only with a hybrid quantum-classical technique based on the generator coordinate method, *Phys. Rev. C* **109**, 024327 (2024).
- [69] M. Q. Hlatshwayo, J. Novak, and E. Litvinova, Quantum benefit of the quantum equation of motion for the strongly coupled many-body problem, *Phys. Rev. C* **109**, 014306 (2024).
- [70] S. Barison, J. R. Moreno, and M. Motta, Quantum-centric computation of molecular excited states with extended sample-based quantum diagonalization, *Quantum Sci. Technol.* **10**, 025034 (2025).
- [71] D. Marti-Dafcik, H. G. A. Burton, and D. P. Tew, Spin coupling is all you need: Encoding strong electron correlation in molecules on quantum computers, *Phys. Rev. Res.* **7**, 013191 (2025).
- [72] H. R. Grimsley and F. A. Evangelista, Challenging excited states from adaptive quantum eigensolvers: Subspace expansions vs. state-averaged strategies, *Quantum Sci. Technol.* **10**, 025003 (2025).
- [73] J. Zhang and D. Lacroix, Excited States from ADAPT-VQE convergence path in Many-Body Problems: application to nuclear pairing problem and  $H_4$  molecule dissociation (2025), arXiv:2506.22275 [quant-ph].
- [74] E. Fradking, *Field Theories of Condensed Matter Physics* (Cambridge University Press, 2013).
- [75] B. Anselme Martin, T. Ayrat, F. m. c. Jamet, M. J. Rančić, and P. Simon, Combining matrix product states and noisy quantum computers for quantum simulation, *Phys. Rev. A* **109**, 062437 (2024).
- [76] J. Dukelsky, S. Pittel, and G. Sierra, Colloquium: Exactly solvable Richardson-Gaudin models for many-body quantum systems, *Rev. Mod. Phys.* **76**, 643 (2004).
- [77] J. von Delft and D. Ralph, Spectroscopy of discrete energy levels in ultrasmall metallic grains, *Phys. Rep.* **345**, 61 (2001).
- [78] V. Zelevinsky and A. Volya, Nuclear pairing: New perspectives, *Phys. of At. Nucl.* **66**, 1781 (2003).
- [79] D. M. Brink and R. A. Broglia, *Nuclear Superfluidity Pairing in Finite Systems* (Cambridge University Press, 2005).
- [80] E. A. Ruiz Guzman and D. Lacroix, Accessing ground-state and excited-state energies in a many-body system after symmetry restoration using quantum computers, *Phys. Rev. C* **105**, 024324 (2022).