

# Hybrid Reward-Driven Reinforcement Learning for Efficient Quantum Circuit Synthesis

Sara Giordano\*, Kornikar Sen\*, and Miguel A. Martin-Delgado\*†

\* Departamento de Física Teórica, Universidad Complutense de Madrid, 28040 Madrid, Spain

† CCS-Center for Computational Simulation, Campus de Montegancedo Universidad Politécnica de Madrid (UPM), 28660 Boadilla del Monte, Madrid, Spain

**Abstract**—A reinforcement learning (RL) framework is introduced for the efficient synthesis of quantum circuits that generate specified target quantum states from a fixed initial state, addressing a central challenge in both the NISQ era and future fault-tolerant quantum computing. The approach utilizes tabular Q-learning, based on action sequences, within a discretized quantum state space, to effectively manage the exponential growth of the space dimension. The framework introduces a hybrid reward mechanism, combining a static, domain-informed reward that guides the agent toward the target state with customizable dynamic penalties that discourage inefficient circuit structures such as gate congestion and redundant state revisits. By leveraging sparse matrix representations and state-space discretization, the method enables scalable navigation of high-dimensional environments while minimizing computational overhead. Benchmarking on graph-state preparation tasks for up to seven qubits, we demonstrate that the algorithm consistently discovers minimal-depth circuits with optimized gate counts. Moreover, extending the framework to a universal gate set for arbitrary quantum states still produces minimal depth circuits, highlighting the algorithm’s robustness and adaptability. The results confirm that this RL-driven approach efficiently explores the complex quantum state space and synthesizes near-optimal quantum circuits, providing a resource-efficient foundation for quantum circuit optimization.

## I. INTRODUCTION

The design of optimized quantum circuits is a crucial task in the current NISQ era and for the future of fault-tolerant quantum computing. Various methods based on combinatorial optimization and mathematically rigorous techniques have been employed for the interest in quantum computing over the past decades. Alongside these standard approaches, the use of Artificial Intelligence (AI)—in particular, Reinforcement Learning (RL) techniques [1]–[3]—has gained importance in the field, contributing to the generation of new circuits with reduced depth, fewer gates, or improved geometric structures tailored to current hardware constraints [4], [5].

In this work, we address the problem of synthesizing quantum circuits that generate a target quantum state from a fixed initial state by means of reinforcement learning. Specifically, we propose a tabular Q-learning framework [6] in which the agent learns action sequences over a discretized representation of quantum states. The learning process is guided by a hybrid reward mechanism: we use a sparse static reward precomputed offline, guiding the agent to the target state, and dynamic penalties calculated during the learning process. These penalties encode other circuit constraints aiming to

optimize simultaneously more features of the final circuit. These components allow the agent to progressively discover efficient paths reaching the target state while avoiding circuit structures that would increase the overall depth [7].

Although the quantum state space grows doubly exponentially with the number of qubits and the action space scales polynomially, our use of discretization and sparse matrix representations [8] keeps the problem tractable under certain assumptions. A Q-learning setup, with a fixed initial state and finite action sequences ensures efficient exploration despite the vastness of the state space [9]. Moreover, we introduce multiple strata of static reward around the target state, in order to reduce its sparsity, and thus improving the algorithm performances. To limit offline computation, we restrict our benchmarks to systems of up to 7 qubits.

The paper is organized as follows. In Sec. II, we review the relevant background on quantum state representations, gate sets, and circuit metrics. Sec. III presents the Q-learning algorithm in detail, including the reward design and exploration strategy. Sec. IV outlines the implementation aspects while Subsection IV-B focuses on graph-state preparation tasks, and Subsection IV-L extends the framework to more general target states. Finally, conclusions and future directions are discussed in Sec. V.

## II. PRELIMINARIES

Before going into the technical details of the algorithm used to build efficient circuits, let us first discuss a few basic concepts which will be used in the remaining part of the paper.

### A. States With Equal-amplitude and Encoded-phase Terms

To utilize the RL method for constructing an appropriate circuit that is capable of producing the exact form of an  $n$ -qubit target state, we should be able to label all possible  $n$ -qubit states contained in the  $n$ -qubit Hilbert,  $\mathcal{H}_n$ . Since  $\mathcal{H}_n$  consists of an infinite number of states, it is infeasible to label each of them. Hence, instead of considering the entire Hilbert space,  $\mathcal{H}_n$ , we focus on a specific finite subset  $\mathcal{S}_n$  of  $\mathcal{H}_n$ , which includes all the states  $|\psi_n\rangle$  that can be expressed as

$$|\psi_n\rangle = \frac{1}{\sqrt{N'}} \sum_{j=1}^N \alpha_j |x_j\rangle, \quad (1)$$

where  $|x_j\rangle$  is an element of the computational basis,  $\mathcal{B}_n \equiv \{|0\rangle, |1\rangle\}^{\otimes n}$  of  $\mathcal{H}_n$   $\forall j$ ,  $\alpha_j \in \mathcal{P}_M := \{e^{i2m\pi/M}\}_{m \in \mathbb{N}}$  is

a complex phase,  $N \in \mathbb{N}$  specifies the number of terms in superposition and  $N'$  is the normalization constant. Here,  $M \in \mathbb{N}$  characterizes the size of the finite set,  $\mathcal{P}_M$ , and correspondingly the discreteness of the complex phases. Care must be taken to recognize that the index  $j$  labels the different terms in the expansion of  $|\psi_n\rangle$ . Importantly, the equality  $|x_j\rangle = |x_{j'}\rangle$  does not necessarily imply that  $j = j'$ . Consequently, the total number of unique indices  $N'$  is, in general, not equal to  $N$ .

This construction allows us to represent states with equal amplitude terms and discrete phases, deliberately neglecting the full complex amplitudes. We refer to such states as States With Equal-amplitude and Encoded-phase Terms (SWEET). We always choose the discreteness,  $M$ , of the phase set,  $\mathcal{P}_M$ , in such a way that the cardinality,  $P_M$ , of  $\mathcal{P}_M$  can be written as  $P_M = M = 2^p$  where  $p \in \mathbb{N}$ . Although this limits the expressive power of  $\mathcal{S}_n$ , it is particularly well suited for representing certain classes of quantum states, such as graph states [10], [11], GHZ states [12], W states [13], and other highly entangled states with uniform amplitudes, where the relative phase carries all relevant information.

Most importantly, this simplification enables a compact representation of quantum states and allows for an efficient integration with our reinforcement learning-based circuit synthesis framework. It leads to a finite and discrete state space suitable for encoding in Q-learning. Moreover, it supports the application of a certain universal gate set (e.g., CNOT,  $H$ ,  $T$ ) in an efficient and tractable manner. Thus, in this work, we limit our study to the SWEET set,  $\mathcal{S}_n$ , due to both computational tractability, appropriateness for our Machine Learning (ML) technique, and its relevance for the targeted class of quantum states.

### B. Universal gates

Any quantum circuit applied to qubit systems can always be built using a fixed finite set of gates or unitaries. This set of gates, say  $\mathcal{U}$ , is called universal gates.  $\mathcal{U}$  is not unique and can be modified according to available resources. In the context of quantum advantage, qubit gates can be classified into two groups, viz. Clifford and non-Clifford gates. Gates that map the Pauli group onto itself, i.e., the gates which normalize the Pauli group, are called Clifford gates. In addition to the Pauli gates themselves, Hadamard ( $H$ ), CZ( $c, t$ ), and CNOT( $c, t$ ) are also some examples of Clifford gates. Here

$$H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \quad (2)$$

and CZ( $c, t$ ) (CNOT( $c, t$ )) denotes the application of  $Z$  ( $X$ ) gate on the target qubit,  $t$ , when the control qubit,  $c$ , is in state  $|1\rangle$ , otherwise acting as the identity on  $t$ .

We know from the Gottesman–Knill theorem [14] that circuits involving only Clifford gates can be simulated on classical computers. Therefore, to introduce computational quantumness in a circuit, one must add non-Clifford gates (gates that do not normalize the Pauli group) in the circuit, which suggests the inclusion of at least one non-Clifford gate

in the set  $\mathcal{U}$ . One of the most popular single-qubit non-Clifford gates is

$$T = \begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{pmatrix}. \quad (3)$$

To create entanglement between the qubits, we need at least one gate in the universal set that can produce entanglement, for example, the CNOT gate. Any set of gates that involves a class of gates that is dense in SU(2), one entangling gate, and one non-Clifford gate can be used as a universal set of gates. For building circuits to produce any  $|\psi_n\rangle \in \mathcal{S}_n$ , we will consider  $\mathcal{U} = \{H, \text{CNOT}, T\}$  as the set of universal gates [15], [16]. Therefore, in the rest of the paper, the notation  $\mathcal{U}$  will be used to denote only this particular universal set of gates.

### C. Circuit complexity

The circuit complexity of a quantum algorithm analyzes the difficulty in building and using that algorithm. Circuit complexity plays a crucial role in a wide range of fields [17], [18]. It serves as a useful tool for understanding the dynamics of quantum systems, especially in the context of chaos [19], decoherence [20], and phase transitions [21].

The most traditional way of quantifying circuit complexity is by counting the total number of single- or two-qubit gates used in the circuit. However, circuit complexity can also be formulated in a continuous geometric framework [22], [23], where it is interpreted as the minimal path in the space of unitary operations needed to transform the identity into a given target unitary. Several approaches have been developed to describe this continuous notion of complexity [24]–[26]. One can also define circuit complexity of quantum algorithms concerning specific resources like entanglement or magic. In the following part, we discuss some of the circuit complexity measures that will be useful in analyzing the RL technique under consideration.

- *Gate count:* Gate count is the total number of gates used in the circuit. It is the most basic measure of circuit complexity. However, this measure provides equal weight to every gate, but in reality, particular gates can be more expensive than the others. That is why other measures, which are discussed below, are also useful.
- *T count:* Non-Clifford gates are more expensive for fault-tolerant computation. Hence, the number of non-Clifford gates used in a circuit measures its complexity. In this work, we will use only one specific non-Clifford gate, i.e., the  $T$  gate, and quantify the circuit complexity as the total number of  $T$  gates used in the circuit. We name this particular type of complexity as  $T$  count.
- *Entangling power:* Since entanglement is an essential quantum resource, the total number of entangling gates involved in the circuit also serves as a measure of complexity, which we refer to as entangling power.
- *Circuit depth:* A quantum circuit may consist of gates that can be applied in parallel, specifically in situations where they operate on different parts of the entire system. The simultaneous action of such gates will reduce the

execution time of the circuit, thus reducing the probability of having errors. We can divide the entire circuit into different time steps, where each time step can involve multiple gates that can be applied simultaneously. Circuit depth is defined as the total number of time steps in a circuit.

### III. ACTION SEQUENCE-BASED Q-LEARNING USING $\epsilon$ -GREEDY EXPLORATION AND HYBRID REWARDS

After defining and discussing all the necessary tools, in this section, we describe the reinforcement learning technique that we use to generate efficient quantum circuits.

#### A. Q-learning basics

RL is used in machine learning to train an agent within a certain environment by giving it feedback based on its actions in the environment [2]. The RL technique that we use is called Q-learning, which can be used to teach the agent to find the best sequence of actions to perform in order to reach a specific state in the environment, without having any prior knowledge of the environment features.

More specifically, in Q-learning, the agent has no prior knowledge of the goodness of an action  $a$  performed in a specific state  $s$ , which implies that Q-learning is an off-policy RL algorithm [6]. A policy is usually denoted as  $\pi(s, a)$  and represents the value assigned to the pair action-state, based on the goodness of this choice, with respect to the target or goal. In Q-learning, the policy is usually represented by the Q-values  $Q(s, a)$ . This makes it different from other policy-based RL, which aims to optimize a preexisting policy by interacting with the environment. Thus, the objective is to define the best policy that will give the optimal path leading to a chosen final state from an initial state.

We can optimize different features of this path by choosing reward and penalty criteria and evaluating the values to feed back to the agent when it selects its actions. The reward function can be represented by a matrix, i.e., the R-matrix, which associates actions and states to a positive, negative, or zero value. This reward assignment can be done in an on-line or off-line fashion, by building a reward function  $R(s, a)$ , before the Q-learning starts or during the training process.

The Q-learning procedure is divided into two parts. The first, the training phase, consists of exploring the environment, i.e., choosing a random action,  $a_t$ , and a random state,  $s_t$ , and building the quality matrix (Q-matrix) using the following Q-update rule:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R(s_t, a_t) + \gamma \cdot \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)], \quad (4)$$

where  $s_{t+1}$  denotes the state that the agent reaches by performing  $a_t$  at  $s_t$ ,  $\alpha$  is the learning rate and  $\gamma$  is the discount factor. The Q-matrix finally specifies the goodness of the state-action pairs, thus the policy to follow to reach the desired state with the optimized path.

The second part, the testing phase, consists of following the path indicated by the Q-values of the Q-matrix and checking if

the learning process brings the agent to the desired target and if the obtained path respects the optimization criteria: a starting point is chosen, and the agent builds the chain of actions to perform in order to reach the final position, selecting in the Q-matrix the highest valued actions for each state it crosses. If it manages to reach the final position meeting the desired criteria for the path, then the learning process is considered successful, and the Q-matrix now represents a successful policy for the agent.

#### B. Q-learning for circuit optimization: states and actions

Let us now directly move to the problem at hand and discuss the Q-learning method in our specific circuit optimization context.

The states  $s \in \mathcal{S}$  and actions  $a \in \mathcal{A}$  that our agent evaluates correspond, respectively, to quantum states and quantum gates. We want to focus on the  $n$ -qubit SWEET set as described in Sec. II-A. To represent the discretized phase  $\alpha_j$  of each term  $j$  in the quantum superposition state  $|\psi_n\rangle$ , we introduce an auxiliary register containing  $p$  additional qubits, supplementing the existing  $n$ -qubit system. Within this architecture, the  $n$ -qubit register encodes the computational basis state  $|x_j\rangle$  that appears in the superposition  $|\psi_n\rangle$ . Since the phase  $\alpha_j$  belongs to the set  $\mathcal{P}_M$ , it can be expressed in the form  $\alpha_j = e^{i(2\pi m_j/M)}$ , where  $m_j$  is a natural number satisfying  $0 \leq m_j < M = 2^p$ . The register of  $p$  qubits registers the binary representation of the number  $m$ . For example, if  $M = 8 = 2^3$  then the state  $|00\rangle - |11\rangle$  will be represented by  $|00000\rangle, |10111\rangle\}$  where the 3 leftmost qubits for each term represent the phase  $000 \rightarrow 1$  and  $101 \rightarrow e^{i\pi} = -1$ , while the 2 rightmost qubits represent the real qubits. We would like to emphasize that there are no extra qubits involved in the circuit; we have introduced it only for the simulation of the discretized phases of the quantum states. This method is inspired by the sign-magnitude representation of binary numbers and allows us to implement gates that can alter the phases of the states and correctly store the result of the gate application, such as  $T$  or  $H$ . Finally, we assign each of these SWEET states an index  $s$ . Since we have  $p$  and  $n$  qubits to register  $\alpha_j$  and  $|x_j\rangle$ , respectively, the total number of forms that a particular term  $\alpha_j |x_j\rangle$ , of  $|\psi_n\rangle$  can take is  $2^{n+p}$ . Since we consider that  $j \neq j'$  implies  $\alpha_j |x_j\rangle \neq \alpha_{j'} |x_{j'}\rangle$ , the number of SWEET states having  $N$  number of terms is  $\binom{2^{n+p}}{N}$ . Hence, the cardinality of the set,  $\mathcal{S}_n$ , of all SWEET states representing an  $n$ -qubit system is

$$N_S = \sum_{N=1}^{2^{n+p}} \binom{2^{n+p}}{N} = 2^{2^{n+p}} - 1. \quad (5)$$

Hence, the number of possible values that the index,  $s$ , which indicates each state, can take is  $N_S$ . For example, when  $n = 7$  and  $p = 1$  (as we will consider to create graph states)  $N_S = O(10^{77})$ .

Meanwhile, the set of available actions  $\mathcal{A}$  (see Sec. III-D) is defined as all valid applications of quantum gates from a fixed set over the  $n$  qubits. Let us denote the number of distinct

gates that jointly act on exactly  $k$  qubits among  $n$  qubits by  $g_k^n$ . For example, for  $\mathcal{U} = \{H, T, \text{CNOT}\}$ ,  $g_1^n = 2$  and  $g_2^n = 1$ .

The total number of valid actions associated with  $k$ -qubit gates is given by  $A_k = g_k^n \cdot \text{Perm}(n, k)$ , where  $\text{Perm}(n, k) = \frac{n!}{(n-k)!}$  is the number of ordered  $k$ -tuples of distinct qubits that can be used as gate arguments (since the order matters for most gates, such as controlled gates). Thus, the total number of actions  $N_A$  is:

$$N_A = |\mathcal{A}| = \sum_k g_k^n \cdot \text{Perm}(n, k). \quad (6)$$

Hence, from the cardinality of the sets  $\mathcal{A}$  and  $\mathcal{S}_n$  we realize that the total dimension of the reward and quality matrices is  $N_S \times N_A$ . This expression grows exponentially with the number of qubits and the richness of the gate set. Combined with exponential growth in the state space, this further motivates the need for efficient exploration strategies and sparse Q-value storage.

Therefore, to store rewards and Q-values efficiently, using the least amount of memory, we use Structured Query Language (SQL) databases [27]. These are relational database systems designed to store and query structured data using tables. Each *table* is defined by a set of *columns* (fields) and stores data in *rows*, where each row represents a distinct record or entry.

In our implementation, SQL is used to efficiently represent large, sparse matrices, such as  $R$  and  $Q$ . Each nonzero element of these matrices, say  $R(s, a)$  or  $Q(s, a)$ , is stored as a table of triplets as

$$\{\text{state\_index}, \text{action\_index}, \text{value}\}. \quad (7)$$

This structure allows us to store only non-zero entries, significantly reducing memory usage. Notice that SQL databases are specifically efficient for queries such as reading or updating specific entries, which are the required actions for our algorithm.

Having in mind the initial state,  $|\psi_0\rangle$ , the target state  $|\psi_{\text{target}}\rangle$ , and the list of actions  $\mathcal{A}$ , we are ready to proceed.

### C. Static and dynamic reward

We employ a hybrid reward strategy in which a sparse and static precomputed reward encourages the discovery of the shortest circuit reaching the target state, while a dynamic action-dependent reward function helps optimize secondary circuit properties such as preparation depth [28]. This choice of using offline and online strategies helps us focusing on the environment exploration and avoid superfluous evaluations during online training of the algorithm, and aligns with the multi-objective reinforcement learning framework [29] and the so-called "MaxPain" techniques [30]. Each element  $R(s, a)$ , of the reward matrix, is split as follows:

$$R(s, a) = R_{\text{sta}}(s, a) + R_{\text{dyn}}(s, a). \quad (8)$$

In Eq. (8), the static reward is precomputed and assigned as follows:

$$R_{\text{sta}}(s, a) = \begin{cases} R_{\max} & \text{if } s \xrightarrow{a} s_{\text{target}} \\ \frac{R_{\max}}{2} & \text{if } \exists a, a_1 : s \xrightarrow{a} s_1 \xrightarrow{a_1} s_{\text{target}} \\ \vdots & \vdots \\ \frac{R_{\max}}{2^k} & \text{if } \exists a, a_1 \dots a_{k-1} : s \xrightarrow{a, a_1 \dots a_{k-1}} s_{\text{target}} \\ 0 & \text{otherwise,} \end{cases} \quad (9)$$

where  $R_{\max} > 0$  is the maximum reward assigned and  $s_{\text{target}}$  is the index representing  $|\psi_{\text{target}}\rangle$ . The other static rewards and penalty values are determined based on this reference value [2], [31]. To ensure that the algorithm functions effectively, penalties should not disrupt the main reward-based optimization. Therefore, individual penalties are set at approximately 1% or 0.1% of  $R_{\max}$ . The reward layers help the agent reach the target state by creating a "breadcrumb trail" for it.

To efficiently compute this static layered reward, we utilize the reversibility property of quantum gates [15]. Say  $U$  is a quantum gate operator, then we have:  $UU^\dagger = U^\dagger U = I$ . We apply the Hermitian conjugate of a gate corresponding to column  $j$ , to the target state:  $U_j^\dagger |\psi_{\text{target}}\rangle = |\psi_i\rangle$  and fix the corresponding element of the static part of the R-matrix as  $R_{\text{sta}}(s_i, a_j) = R_{\max}$ , where  $s_i$  and  $a_j$  represent the state  $|\psi_i\rangle$  and the action of  $U_j$  on  $|\psi_i\rangle$ , respectively. To construct the second layer, we apply  $U_k^\dagger$  to  $|\psi_i\rangle$ , resulting in the states for the second stratum  $|\psi_l\rangle$  and set  $R_{\text{sta}}[s_l, a_k] = R_{\max}/2$ . This process is repeated for the remaining strata. By repeating this procedure for all gates, we can build the reward matrix successfully. Here we write the steps of the reward assignment:

1. Define the target state  $|\psi\rangle_{\text{target}}$ , the maximum reward value as  $R_{\max} = 10000$  and the number of strata  $k_{\max}$ .
2. Take  $U_j$  from the available set of gates and compute  $U_j^\dagger$ .
3. Apply  $U_j^\dagger$  to the current target state  $|\psi_{\text{target}}\rangle$ , obtaining a new state  $|\psi_i\rangle$ .
4. Let the current recursion strata be  $k$ . If  $R_{\text{sta}}(s_i, a_j) < R_{\max}/2^k$ , set  $R_{\text{sta}}(s_i, a_j) = R_{\max}/2^k$ .
5. While  $k < k_{\max}$ , apply steps 2–5 recursively with the updated state  $|\psi_{\text{target}}\rangle \leftarrow |\psi_i\rangle$ ,  $k \leftarrow k+1$ , and all possible gates  $U_j$ .
6. Repeat for all possible gates  $U_j$ .

After labeling each state,  $|\psi_i\rangle \in \mathcal{S}_n$  with an index  $s_i$ , and each gate,  $U_j$ , acting on fixed qubits with an action  $a_j$ , we can assign rewards using the procedure outlined in Algorithm 1. The dynamic component in Eq. (8) is computed during Q-learning training, based on the sequence of actions taken in each episode and their varying impacts on the environment. It involves penalties assigned to the agent for following ineffective paths. Specifically, penalties are applied for:

- i. revisiting states within the same episode,
- ii. performing actions that do not change the state,

In addition, we can incorporate other penalties based on specific requirements. For instance, penalties can be added to discourage the use of costly gates or to reduce the circuit

---

**Algorithm 1** Recursive assignment of static reward values

```

1: Set  $s_{\text{target}}$ 
2: Set max reward  $R_{\max} \leftarrow 10000$ 
3: Set maximum number of strata  $k_{\max}$ 
4: procedure ASSIGNSTATICREWARD( $s_{\text{target}}, k$ )
5:   for each action  $U_j$  in the available gate set do
6:     Compute inverse  $U_j^\dagger$ 
7:     Apply  $a_j$  corresponding to  $U_j^\dagger$  to  $s_{\text{target}}$ , and get  $s_i$ 
8:     if  $R_{\text{sta}}(s_i, a_j) < R_{\max}/2^k$  then
9:       Set  $R_{\text{sta}}(s_i, a_j) \leftarrow R_{\max}/2^k$ 
10:    end if
11:    if  $k < k_{\max}$  then
12:       $k \leftarrow k + 1$ 
13:       $s_{\text{target}} \leftarrow s_i$ 
14:      ASSIGNSTATICREWARD( $s_{\text{target}}, k$ )
15:    end if
16:  end for
17: end procedure

```

---

depth. These specific examples are discussed in the next section.

During each training episode, to impose penalties for i. revisiting states and ii. taking ineffective actions, we record the visited states in a list called  $V$ . We check  $V$  to identify revisited or unchanged states (see Eq.(10)). These penalties are directly incorporated into the reward signal and included in the standard Q-learning update.

$$R_{\text{dyn}}(s_t, a_t) = \begin{cases} R_{\text{dyn}}(s_t, a_t) - R_{\max} \cdot 10^{-4} & \text{if } s_t \xrightarrow{a_t} s_{t+1} \\ & \text{where } s_{t+1} \in V \\ & \text{and } R_{\text{sta}}(s_t, a_t) = 0 \\ R_{\text{dyn}}(s_t, a_t) - R_{\max} \cdot 10^{-3} & \text{if } s_t \xrightarrow{a_t} s_{t+1} = s_t. \end{cases} \quad (10)$$

After outlining the process of constructing the reward matrix, we now move on to the details of the agent training method.

#### D. Training phase for action-sequence based Q-learning

Our approach utilizes the standard tabular Q-learning algorithm, which iteratively updates the value function  $Q(s, a)$  through temporal difference (TD) learning and the traditional Q-update rule. This is achieved by randomly selecting gates to apply to randomly chosen states and exploring the entire state-action space. [1], [2].

The training process is organized into episodes. Each episode begins with the agent starting from a fixed initial state  $s_0$ , representing a specific quantum state  $|\psi_0\rangle \in \mathcal{S}_n$ . The agent then performs a sequence of  $k$  actions, updating the state after each action. Notably, in this setup, the initial state is not randomly chosen, as in standard Q-learning. At each step, the agent can select an action either randomly from the set of possible actions  $\mathcal{A}$  (exploration) or by choosing the action with the highest current Q-value (exploitation). To balance exploration and exploitation, we use the  $\epsilon$ -greedy strategy [2]: with probability  $\epsilon$ , the agent explores by selecting a random

---

**Algorithm 2** Training phase of the Q-learning algorithm

```

1: Define learning parameters: exploration rate  $\epsilon$ , learning rate  $\alpha$ , discount factor  $\gamma$ 
2: Initialize the Q-matrix
3: Set number of episodes  $l_e$  and episode length  $k$ 
4: Define the set of possible actions  $\mathcal{A}$ 
5: for each episode  $e = 1$  to  $l_e$  do
   ▷ Fixed initial SWEET state
6:   Initialize  $s_0$  corresponding to  $|\psi_0\rangle$ 
7:   for each step  $t = 0$  to  $k - 1$  do
8:     Select action  $a_t \in \mathcal{A}$  using  $\epsilon$ -greedy strategy:
9:     with probability  $\epsilon$ , choose  $a_t$  randomly
   ▷ Exploration
10:    with probability  $1 - \epsilon$ , choose  $a_t = \max_a Q(s_t, a)$ 
   ▷ Exploitation
11:    Apply  $a_t$  to  $s_t$ , obtaining new state  $s_{t+1}$ 
12:    Compute  $R(s_t, s_{t+1}) = R_{\text{sta}}(s_t, a_t) + R_{\text{dyn}}(s_t, a_t)$ 
13:    Update Q-value using temporal difference:
      
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R(a_t, s_t) +$$

      
$$+ \gamma \cdot \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

14:    Set  $s_t \leftarrow s_{t+1}$ 
15:  end for
16: end for
17: Proceed to testing phase (see Algorithm 3)

```

---

action, and with probability  $1 - \epsilon$ , it exploits by choosing the best-valued action. This sequential process maintains the off-policy nature of Q-learning while promoting exploration and learning in high-dimensional environments.

We denote the sequence of actions in one episode by  $A$ . After completing  $l_e$  episodes of fixed length  $k$ , the algorithm transitions to the testing phase. The values  $k$  and  $l_e$  depend on the specific application. Training continues with batches of episodes until the testing phase proves successful (see Sec. III-E). The primary objective of training is to iteratively build and update the Q-matrix,  $Q$ , which represents the expected long-term reward for taking action  $a$  in state  $s$ . During each episode, after every applied action, the agent receives a reward signal composed of both static and dynamic components, as detailed previously in Eq. (4) and Sec. III-C. The static reward, computed offline, directs the agent toward the target state by rewarding minimal length paths and is stored as a sparse matrix. The dynamic reward is calculated online during training, with the elements of the matrix  $R_{\text{dyn}}(s, a)$  constructed according to Eq. (10).

These reward signals are used to update the Q-values using the temporal difference learning rule in Eq. (4), where  $R(s_t, a_t)$  represents the total reward obtained at step  $t$ , as described in Eq. (8). Across multiple episodes, the Q-matrix progressively encodes the most efficient action sequences for reaching the target quantum state, optimizing not only for gate count but also for any additional criteria if specified.

Our training structure aligns with Q-learning based on

action sequences and planning-inspired exploration, similar to models like Dyna-Q [32]–[34]. However, our algorithm remains strictly off-policy and tabular: it does not utilize neural network approximators as in Deep Q-learning [35], nor does it directly optimize a policy as in policy gradient or actor-critic methods [36], [37]. This setup enables more efficient exploration in high-dimensional spaces while still allowing Q-learning to converge to an optimal solution under suitable exploration strategies.

The training phase occurs before the testing phase (see Sec. III-E), during which the effectiveness of the algorithm is evaluated. The training process is repeated in batches until satisfactory performance is achieved in the testing phase. This sequential approach, together with the use of hybrid rewards and a discretized representation of quantum states, facilitates effective learning in high-dimensional and sparse environments.

#### E. Q-learning testing phase

To evaluate the learned policy matrix  $Q$ , we generate a circuit based on  $Q$  after completing a batch of episodes and then assess the training success. The testing phase begins by selecting the row  $s_0$ , representing the initial state  $|\psi_0\rangle$ , and setting a maximum circuit length  $l_c$ . We identify the action  $a_0$  for which  $Q(s_0, a_0) \geq Q(s_0, a)$  for all  $a$ . The corresponding gate  $U_0$ , represented by  $a_0$ , is applied to  $|\psi_0\rangle$ , resulting in  $|\psi_1\rangle = U_0 |\psi_0\rangle$ . We then check if  $|\psi_1\rangle$  matches the target state  $|\psi_{\text{target}}\rangle$ . If not, we find the row  $s_1$  corresponding to  $|\psi_1\rangle$  and repeat the process: selecting the action  $a_t$  that maximizes  $Q(s_t, a)$ , applying the associated gate  $U_t$ , and updating the state to  $|\psi_{t+1}\rangle$ . This cycle continues until the state matches the target  $|\psi_{\text{target}}\rangle$ , or until the maximum circuit length  $l_c$  is reached (i.e.,  $t \leq p$ ). If successful, the algorithm concludes; otherwise, the process returns to training.

Hence, the testing phase is as follows.

1. Find the row  $s_i$ , which represents the state  $|\psi_i\rangle$ .
2. Find the column  $a_i$ , for which  $Q(s_i, a_i) \geq Q(s_i, a)$  for all  $a$ .
3. Apply the gate,  $U_i$ , represented by the column,  $a_i$ , on  $|\psi_i\rangle$ .
4. Check if  $U |\psi_i\rangle$  equals the target state. If it does, the testing phase terminates. If not and the total circuit length is less than  $l_c$ , reset  $|\psi_i\rangle$  to  $U_i |\psi_i\rangle$  and return to Step 1.
5. If the circuit length reaches  $l_c$ , the process switches back to the training phase, initiating a new training batch. The Q-matrix used for testing is used to initialize the Q-matrix for this new training batch.

After the algorithm concludes, the sequence of gates corresponding to the actions  $\{a_0, \dots, a_{t-1}\}$  constitutes the final circuit. This circuit transforms the initial state  $|\psi_0\rangle$  into the target state  $|\psi_{\text{target}}\rangle$  using the minimum number of gates, which is less than  $l_c$ , the maximum circuit length specified at the start.

This Q-learning approach not only constructs a circuit with an optimal length but also allows the incorporation of additional constraints through suitable penalties. For instance, we can discourage the use of T-gates, entangling gates, or any

---

#### Algorithm 3 Testing phase of the Q-learning algorithm

---

```

1: Set initial state  $s_0$  corresponding to  $|\psi_0\rangle$ 
2: Set target state  $s_{\text{target}}$  corresponding to  $|\psi_{\text{target}}\rangle$ 
3: Set maximum allowed circuit length  $l_c$ 
4: Initialize step counter  $t \leftarrow 0$ 
5: while  $s_t \neq s_{\text{target}}$  and  $t < k$  do
6:   Select action  $a_t = \max_a Q(s_t, a)$ 
7:   Apply gate corresponding to  $a_t$  to  $s_t$ , and obtain  $s_{t+1}$ 
8:   Set  $s_t \leftarrow s_{t+1}$ 
9: end while
10: if  $s_t = s_{\text{target}}$  then
11:   Sequence  $\{a_0, a_1, \dots, a_{t-1}\}$  defines the final circuit
12: else
13:   Testing failed: maximum circuit length  $l_c$  reached:
      return to training phase in Algorithm 2
14: end if

```

---

specific expensive gate by updating the dynamic component of the reward matrix with penalties. Additionally, appropriate penalties can be employed to minimize circuit depth. These optimization strategies will be discussed in more detail in Sec. IV.

#### F. Exploration complexity

The exploration space generated by our formulation is naturally large because of the combinatorial complexity of the quantum states involved. Given that  $n$  and  $p$  denote the number of qubits and phase qubits used to encode the states, respectively, the full Hilbert space is defined over  $n+p$  qubits. This results in a computational basis  $\mathcal{B}_{n+p}$  of size  $2^{n+p}$ . The total number of states,  $N_S$ , that the agent explores, i.e. the cardinality of  $\mathcal{S}_n$ , is even greater.

Within this vast space, the reward is highly sparse, making exploring the environment particularly challenging. To address this, we do not explicitly store or explore the entire state-action space. Instead, we utilize sparse representations of the Q-matrix and the R-matrix in database-style structures of the form  $\{\text{state index, gate index, value}\}$ , where each entry represents a valid (nonzero) Q-value or reward, as described in Sec. III-B. This approach allows us to efficiently track only the explored or relevant transitions, significantly reducing memory consumption and computational overhead.

Furthermore, employing sequence-based Q-learning with a fixed initial state and structured penalties promotes efficient exploration and helps the agent to reliably converge to the target state [6]. The dynamic penalties—designed to discourage unnecessary actions—stop the agent from aimless wandering and speed up the learning process.

However, the method is still affected by the curse of dimensionality [9]. Although representing the Q-table and the reward matrix as sparse matrices greatly reduces memory requirements and improves storage scalability [8], the sparsity of static rewards presents a different challenge. As the number of qubits grows, the likelihood of reaching a rewarding state-action pair, where  $R(s, a) > 0$ , decreases quickly. To address

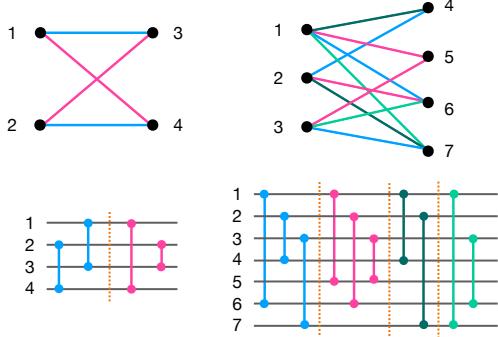


Fig. 1. The figure illustrates the structure of two specific graphs alongside diagrams of the circuits used to generate them. The top panel displays a 4-vertex graph (left) and a 7-vertex graph (right). Edges of the same color signify that the corresponding CZ gates can be applied simultaneously at identical time steps to create the graph states. The bottom panel shows the optimal circuits with minimal depths required to generate these graph states. In the circuit diagrams, vertical solid lines represent the application of CZ gates, while colored circles indicate the qubits involved in each gate's application. Vertical dashed lines separate different time steps.

this, we implement static reward strata at progressively greater distances from the target state.

To address this limitation and keep training computationally manageable, the examples presented in the following sections are limited to systems with up to  $n = 7$ ,  $p = 1$  qubits, and  $n = 5$ ,  $p = 3$ .

#### IV. IMPLEMENTATION AND RESULTS

After reviewing the algorithm details, this section demonstrates how the proposed tabular Q-learning framework converges toward high-quality solutions and how this convergence materializes in concrete quantum circuits. We proceed to illustrate first a benchmarking of the algorithm on graph states, and second a more general result employing a universal set of gates.

##### A. Benchmarking results on graph states

We notice that the chosen SWEET state set is well suited for representing graph states [11]. Graph states form a key class of highly entangled quantum states, widely utilized as resources in measurement-based quantum computation [38]. Their structured entanglement also supports applications in quantum error correction, communication, and multipartite entanglement research [39]. Before analyzing the performance of the learning algorithm, we first clarify the formal structure of graph states and discuss the theoretical background, which makes them ideal benchmarking candidates.

##### B. Graph states definition

A graph,  $G = (\mathcal{V}, \mathcal{E})$ , is defined through a set  $\mathcal{V}$  of vertices (or nodes) and a set of edges,  $\mathcal{E} \equiv \{\{v_i, v_j\} | v_i, v_j \in \mathcal{V}\}$ . For every graph having  $n$  vertices, we can always define a corresponding  $n$ -qubit state as:

$$|G_n\rangle = \prod_{\{v_i, v_j\} \in \mathcal{E}} \text{CZ}(v_i, v_j) |+\rangle^{\otimes n}, \quad (11)$$

hence, called graph state. For an  $n$ -qubit computational basis state  $|x_1 x_2 \dots x_n\rangle$  with  $x_k \in \{0, 1\}$ , the action of the CZ gate on qubits  $i$  and  $j$  is defined as:

$$\text{CZ}(i, j) |x_1 x_2 \dots x_n\rangle = (-1)^{x_i x_j} |x_1 x_2 \dots x_n\rangle. \quad (12)$$

That is, the state acquires a phase of  $-1$  only if both  $x_i = 1$  and  $x_j = 1$ . Thus,  $\text{CZ}(v_i, v_j)$  is the gate CZ acting on the qubit corresponding to vertex  $v_j$  and controlled by the qubit corresponding to vertex  $v_i$ . The initial state used for graph state preparation is  $|+\rangle^{\otimes n}$ , a uniform superposition of the computational basis, and we denote it as  $|\psi_n^{in}\rangle$ . Therefore, the number of CZ gates required to create a graph state from  $|\psi_n^{in}\rangle$  is known and equal to the total number of edges of the graph. CZ gates entangle qubits according to the edges of  $G$ . Since CZ gates commute, the total number of CZ gates equals  $|\mathcal{E}|$ , and the execution order can be rearranged to minimize circuit depth. Graph states can be seen as entangled stabilizer states whose structure is fully determined by the connectivity of the underlying graph.

##### C. Optimal circuit depth for graph states

We now describe the specific target states used to benchmark the reinforcement learning algorithm, along with their known theoretical depth bounds. We will focus on the graph states corresponding to the graphs depicted in Fig. 1,  $|G_4\rangle$  and  $|G_7\rangle$ .

Let us define as  $\mu(|\psi\rangle, |\phi\rangle)$  the number of gates composing a circuit needed to generate state  $|\phi\rangle$  starting from state  $|\psi\rangle$ . To create the graph states  $|G_4\rangle$  and  $|G_7\rangle$  that have 4 and 7 vertices, and 4 and 10 edges, respectively, we need  $\mu(|\psi_4^{in}\rangle, |G_4\rangle) = 4$  and  $\mu(|\psi_7^{in}\rangle, |G_7\rangle) = 10$  CZ gates.

Let us also define, analogously to  $\mu(|\psi\rangle, |\phi\rangle)$ , the depth of the circuit generating  $|\phi\rangle$  from  $|\psi\rangle$  as  $\delta(|\psi\rangle, |\phi\rangle)$ . Using the Vizing theorem [40], it can be shown that the minimum depth,  $\delta(|\psi_n^{in}\rangle, |G\rangle)$ , of the circuit that generates a graph state  $|G\rangle$  using CZ gates, starting from  $|\psi_n^{in}\rangle$ , is directly related to the maximum degree,  $\Delta(G)$ , of the graph,  $G$ . The degree of a vertex is the number of edges connected to the vertex. In particular,  $\delta(|\psi_n^i\rangle, |G\rangle)$  is either  $\Delta(G)$  or  $\Delta(G) + 1$ . For bipartite graphs (like those depicted in Fig. 1) the bound tightens to  $\Delta(G)$ . Hence:

$$\begin{aligned} \delta(|\psi_4^{in}\rangle, |G_4\rangle) &= 2, \\ \delta(|\psi_7^{in}\rangle, |G_7\rangle) &= 4. \end{aligned}$$

These known values will allow us to evaluate whether the learning agent can match the theoretical optimal circuits in terms of both depth and gate count, in particular for states  $|G_4\rangle$  and  $|G_7\rangle$ . The initial states for the circuit generation are  $|\psi_4^{in}\rangle$  and  $|\psi_7^{in}\rangle$ , respectively. The SWEET sets that we focus on are  $\mathcal{S}_4$  and  $\mathcal{S}_7$  (defined in Subsection II-A) and the actions considered,  $\mathcal{A}$ , correspond to CZ gates. We specify that the phase set we are considering is  $\mathcal{P}_2 = \{1, -1\}$ , i.e., we take  $M = 2$  in the SWEET states definition in Eq. (1).  $N_S$  and  $N_A$  can be obtained using Eqs. (5) and (6) by putting  $n = 4$  and 7 respectively for  $|G_4\rangle$  and  $|G_7\rangle$ , and considering  $g_1^n = 0$  and  $g_2^n = 1$ , as defined in Sec. III-F.

#### D. Reward strata and penalty design

The RL setup described in Sec. III is tailored here to favor low-depth circuit discovery for graph states. Specifically, to minimize circuit depth, we define a *congestion penalty* that penalizes actions that reuse qubits already involved in recent gates.

We evaluate how much the qubits involved in an action  $a_t$  have been involved in the past actions of the same episode sequence, defining the “congestion level”  $C(a_t)$  as follows. First, we define a counter vector  $\mathbf{c} = (c_1, \dots, c_n)$ , initialized at the beginning of each episode as  $c_i = 0$  for all  $i \in \{1, \dots, n\}$ . The size of the vector is carefully chosen to be equal to the number of qubits  $n$  involved in the target state. We iterate over all the actions prior to action  $a_t$  in the episode, i.e.  $A_{<t} = \{a_0, \dots, a_{t-1}\}$ . For each qubit  $q_i$  used by  $a_j \in A_{<t}$  we increment the corresponding counters  $c_i \leftarrow c_i + 1$ . Considering that the current action  $a_t$  acts on a subset  $\mathcal{Q}_{a_t}$  of qubits, the congestion level is:

$$C(a_t) = \max_{q_i \in \mathcal{Q}_{a_t}} c_i. \quad (13)$$

We then apply the penalty:

$$R_{\text{dyn}}(s_t, a_t) = R_{\text{dyn}}(s_t, a_t) - R_{\text{max}} \cdot 10^{-4} \quad \text{if } C(a_t) > CT_t, \quad (14)$$

where  $CT_t = t/2$  is a congestion threshold depending on the time step  $t$ . This mechanism discourages serialized actions and encourages parallel gate application.

#### E. Training phase and visualization of the learning process

In this subsection, we analyze the learning dynamic summarized in Fig. 2 which brings to the optimal circuits of Fig. 1.

In particular, Fig. 2 condenses the entire training history for creating  $|G_4\rangle$  into three sparse matrices, the static reward matrix,  $R_{\text{sta}}(s, a)$  (left panel), the learned Q-values  $Q(s, a)$  (centre panel), and dynamic penalty map  $R_{\text{dyn}}(s, a)$  (right panel). Notice that, to facilitate interpretation, the plot includes only state-action pairs with non-negative values.

The left panel encodes the layered breadcrumb trail introduced in Sec. III-C. Only state-action pairs included in the two reward strata described in Eq. (9) have nonzero entries, making the figure visibly sparse. The tallest bars ( $R_{\text{max}} = 10^4$ ) correspond to single-gate transitions that reach the target state directly, while the shorter bars at  $R_{\text{max}}/2$  correspond to states that are two action-steps away. For both states  $|G_4\rangle$  and  $|G_7\rangle$  we set  $R_{\text{max}} = 10^4$  and  $k_{\text{max}} = 2$ , to have two reward strata.

The central panel shows nonzero Q bars highlighting the spread reward with smoothly varying amplitudes. This amplification results from temporal difference updates (see Eq. (4)). We considered the following values of the training parameters both for  $|G_4\rangle$  and  $|G_7\rangle$ :  $\epsilon = 0.8$ ,  $\alpha = 0.8$ ,  $\gamma = 0.5$  for the  $\epsilon$ -greedy strategy and the Q-learning update, respectively. The algorithm provided the circuits in Fig. 1 within a training of  $10^4$  episodes for  $|G_4\rangle$  and  $7 \times 10^4$  for  $|G_7\rangle$ , where each episode involves action sequences of length 50 and starts from the initial state  $|\psi_4^{in}\rangle$  and  $|\psi_7^{in}\rangle$ , respectively.

The right panel shows  $R_{\text{dyn}}$ , containing exclusively negative values due to penalties applied for revisiting states, actions leaving the state unchanged, or exceeding congestion thresholds. The distribution of cool-colored bars confirms that exploration is steered away from inefficient behaviors.

Taken together, the panels confirm convergence: the Q-table inherits the sparse structure of  $R_{\text{sta}}$  and avoids high-penalty regions in  $R_{\text{dyn}}$ . The training hyperparameters were empirically tuned to balance exploration and convergence stability.

#### F. Consideration on average Bellman error

In Q-learning, the Bellman error quantifies the discrepancy between the current estimate of the Q-value for a state-action pair and the updated value predicted by the Bellman equation shown in Sec. III-A. Minimizing this error is fundamental for the algorithm to converge to the optimal Q-table and produce accurate value estimates, and is the result expected after a successful training phase [2], [41], [42]. Specifically, the Bellman error at each update step is given by the temporal difference term:

$$\delta_t = R(s_t, a_t) + \gamma \cdot \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t). \quad (15)$$

By iteratively minimizing  $\delta_t$ , the algorithm refines its Q-values towards the true expected returns, thus learning an optimal policy.

At the end of the training phase of our algorithm, we measure an average Bellman error well below 1% of the maximum reward scale  $R_{\text{max}} = 10^4$ , typically below 100 in absolute value. It is important to stress that this 1% threshold is a practical heuristic rather than a universal theoretical criterion [41], [42]. Such a low average Bellman error confirms that the learned Q-values closely approximate the expected future rewards, signaling effective convergence of the Q-learning process. This precision level aligns with common practical heuristics in reinforcement learning. Consequently, the small Bellman error reinforces confidence in the robustness and near-optimality of the quantum circuits synthesized by the learned policy.

#### G. Testing phase and circuit reconstruction

To validate the learned policy, we run the testing phase (Algorithm 3). Starting with the same initial state of the episodes, the action with the maximum Q-value is selected until the target state is reached or the circuit length cap  $l_c = 50$  is reached.

Fig. 1 shows the final circuits obtained for the two benchmark graphs:

- Four-qubit square graph ( $|G_4\rangle$ ): A depth-2 circuit with 4 CZ gates.
- Seven-qubit bipartite graph ( $|G_7\rangle$ ): A depth-4 circuit with 10 CZ gates.

These results match the optimal theoretical bounds derived from Vizing theorem [40]:  $\delta(|\psi_4^{in}\rangle, |G_4\rangle) = \Delta(G_4) = 2$  and  $\delta(|\psi_7^{in}\rangle, |G_7\rangle) = \Delta(G_7) = 4$ . Thus, the Q-table, shaped by

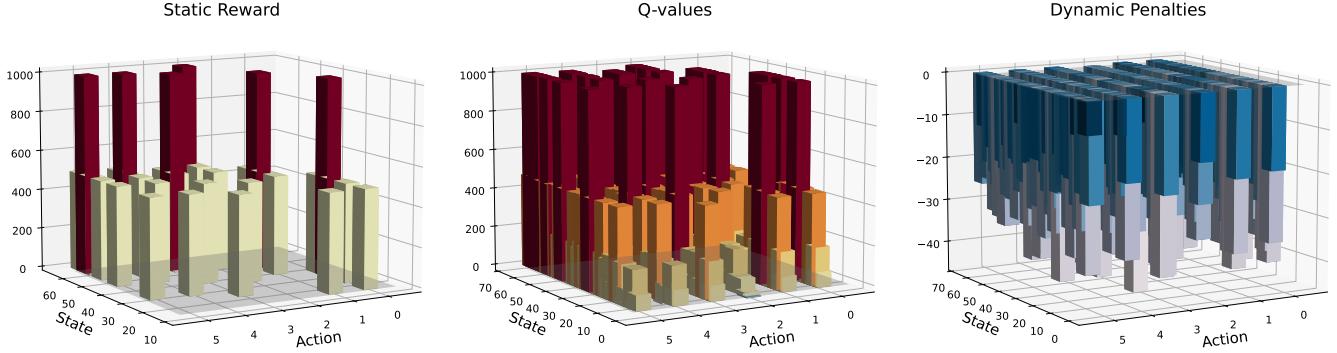


Fig. 2. The learned Q-matrix and its constituent reward components. The figure presents 3D bar charts of the matrices used in the Q-learning algorithm for a 4-qubit graph state. From left to right, the first chart depicts the static reward matrix  $R_{\text{sta}}(s, a)$  with two reward strata. The second shows the learned Q-values  $Q(s, a)$  after training. The third illustrates the dynamic penalty matrix  $R_{\text{dyn}}(s, a)$ , computed online. Color gradients indicate the magnitude and sign of values: warm colors for positive and cool colors for negative. Collectively, these plots illustrate how static guidance and dynamic penalties interact to construct and learn a policy. For more details, see Sec. III-C and Subsection IV-E. For clarity of visualization, only the state-action pairs with non-negative values are shown in the plot.

static rewards and dynamic penalties, encodes the parallelizable structure of these circuits. The learned policy thus yields not only gate-optimal but also discovers the most parallelized structure permitted by the target state, thus attaining depth-optimal results.

#### H. Scalability remarks for graph states

While the results presented involve systems with  $n + p \leq 8$  qubits, where  $n \leq 7$  and  $p = 1$ , the sparse Q and R matrices remain tractable beyond this range. For  $n = 7$  and  $p = 1$ , the combined storage of nonzero entries fits within 170MB. Preliminary results for  $n = 8$  suggest that the method remains effective and efficient due to its reward-guided exploration.

#### I. General results using universal gate set

This RL technique can be used to create any SWEET state and can be customized to optimize different circuit features. The purpose of this section is indeed not only to show the algorithm's ability to target states other than graph states, but also to demonstrate its flexibility in incorporating various constraints through adaptable dynamical penalties, depending on the requirements.

In order to generalize the algorithm, our list of actions will consist in gates from the universal gates set  $\mathcal{U} = \{H, \text{CNOT}, T, T^\dagger\}$  [15], [16]. Notice that, being aware of the equivalence  $T^7 = T^\dagger$ , we intentionally introduce the redundant gate  $T^\dagger$  into the universal set to simplify circuit optimization. To implement the operation of  $T$  and  $T^\dagger$  we need to represent in the SWEET states  $M = 2^3$  different phases (see Eq.(1)), thus we include three phase qubits ( $p = 3$ ) [43], [44]. We choose to target a state with  $n = 3$  qubits such that the set of all possible SWEET states has a dimension  $N_S$  on the order of  $10^{19}$ . The target state is purposely selected to require  $T$  or  $T^\dagger$  gates in its generating circuit, starting from the initial

state  $|\psi_3^{\text{in}}\rangle = |000\rangle$ , in order to target a nontrivial circuit. The chosen target state is as follows:

$$|\psi_3\rangle = |010\rangle + |011\rangle + |100\rangle. \quad (16)$$

One can notice that both  $|\psi_3^{\text{in}}\rangle$  and  $|\psi_3\rangle$  belong to the SWEET states set.

#### J. Designing the reward

In order to highlight the capability of the algorithm to be customized for different optimization objectives, we first consider generating two circuits:

- 1) a circuit involving the minimum number of gates and
- 2) a circuit with optimized depth.

For the first case, we build the static reward matrix with the reward strata discussed in III-C, while, in the dynamic reward, we only include the basic penalties described in Eq.(10). For the second case, we include the depth penalty along with the basic penalties, with the same method used for graph states (see Eqs. (13) and (14)). We consider  $R_{\text{max}} = 10^4$  and four reward strata, i.e.,  $k_{\text{max}} = 4$ .

#### K. Training phase and sparseness of the matrices

We recall that the total number of qubits, including phase qubits, is  $n + p = 6$  and thus  $N_S$  is of the order of  $10^{19}$  and the number of actions is  $N_A = 15$  (see Eqs.(5) and (6)). Despite the vastness of the state-action space, the algorithm is able to construct efficient circuits by exploring only a relevant subspace—sufficient to achieve the objective—while effectively avoiding inconclusive or inefficient paths. This can be seen in Fig. 3, where we depict a projection plot of the R and Q matrices generated at the end of the training phase, for the first type of circuit design, only including basic dynamic penalties. The plot presents a linear scale for the states on the vertical axes, for which the tick labels represent the order of magnitude

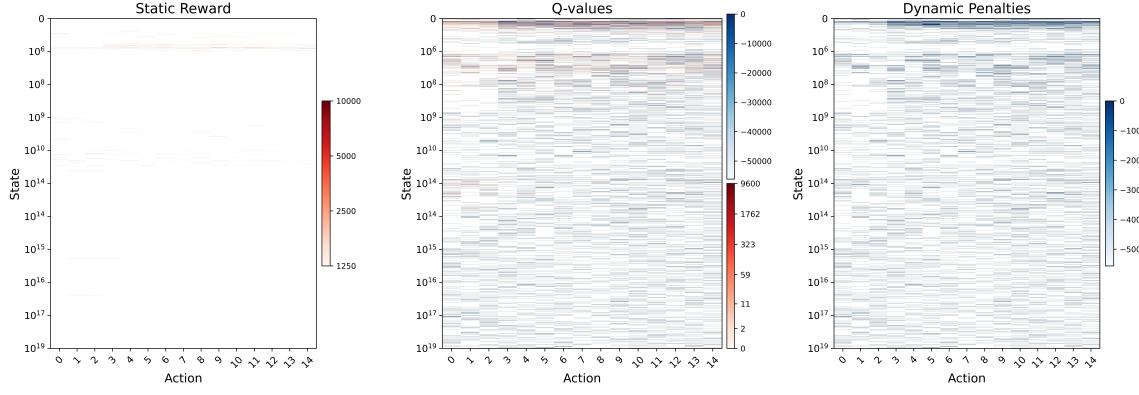


Fig. 3. Visualization of the core Q-learning matrices after training on the state  $|\psi_3\rangle$  (see Eq.(16)). This figure is analogous to Fig.2, but here the much larger state-action space makes its structure more apparent. The left panel shows the sparse static reward matrix  $R_{\text{sta}}$ , which assigns positive "breadcrumb trail" rewards to actions that lead toward the target state. The central panel displays the final Q-matrix after training. Positive values correspond to preferred actions, while negative values indicate actions to be avoided. The right panel presents the dynamic penalty matrix  $R_{\text{dyn}}$ , which captures the negative rewards used during training to discourage inefficient behavior. Together, the static and dynamic reward components shape the Q-matrix, guiding the agent toward an efficient policy. Warm colors denote positive values, and cool colors denote negative values. For further details, refer to Sec.IV-K.

of the index  $s$ . On the horizontal axes are represented the actions indexes. This plot is analogous to the plot in Fig.2, but shows the result of the training phase for a much larger space of  $10^{20}$  state-action pairs, showing both zero and nonzero values of the sparse matrices.

The left panel of the plot represents the static reward matrix  $R_{\text{sta}}$ , which is evidently extremely sparse. It includes only positive values, represented using warm colors. The right panel illustrates the nonzero values of the dynamical reward,  $R_{\text{dyn}}$ , i.e. the penalties applied dynamically during the training. It only includes negative values, represented using cool colors; the brighter the color, the more penalized the agent would be for choosing the corresponding state-action pair. Although the dynamical rewards are spread all over the space, the matrix is again extremely sparse considering the magnitude of the vertical axis values. The central panel shows the Q-matrix. It consists mostly of negative values along with a few positive ones, constituting the agent policy which guides towards the target state. The static rewards contributed to the positive values of the Q-matrix, whereas the dynamical reward provided negative values, according to Eq.(4). We recall that, unlike the plot in Fig. 2, this one displays the three matrices on their original scale; therefore, positive values appear in the Q-matrix, in the same region of the R-matrix where  $R_{\text{sta}}(s, a)$  is also positive. The sparsity of the Q-matrix, clearly visible in the plot, indicates that the agent did not explore the entire space during training but rather focused on specific regions, which led to the fast convergence of the algorithm.

#### L. Optimized circuits with the universal set of gates

The designs of the final circuits determined during the testing phase are illustrated in Fig.4. The top panel shows the circuit found using the basic penalties defined in Eq.(10). For this case, we trained the agent over  $30 \times 10^3$  episodes of length 50, thus around  $15 \times 10^5$  training steps. One can

notice that the circuit built through the algorithm includes only 13 gates; however, the depth of the circuit is 11. To optimize the depth, we included the depth penalties described in Subsection IV-D, and constructed the circuit shown in the bottom panel of the same figure. Since the first circuit required only 13 gates, for the latter case, we used episodes of shorter length, that is, of length 30 instead of 50. Formulating an efficient circuit required training over  $72 \times 10^3$  episodes, thus  $21.6 \times 10^5$  training steps, comparable to the training required with the basic penalties.

As can be noticed from Fig.4, the inclusion of the penalty successfully reduced the depth of the circuit from 11 to 7, at the cost of incorporating more gates, which increased from 13 to 15. Hence, these examples further prove the efficiency of the algorithm in modeling and shaping the circuits according to the user's requirements.

To avoid misinterpretation of the results, it is important to clarify the relationship between the final quantum state produced by the optimized circuit and the targeted SWEET state. During training, when the agent encounters non-SWEET states, it maps them onto the discretized SWEET states, an essential step for efficient exploration of the Hilbert space. Such non-SWEET states may arise, for instance, through the application of non-trivial gates such as Hadamard gates. As a result, the final state generated by the circuit may not be exactly equal to the originally targeted SWEET state, but is associated with it via the discretization process. In this sense, the targeted SWEET state serves as the representative of the actual output state produced by the optimized circuit.

The optimal circuits shown in Fig. 4, which were generated using the Q-learning technique and applied to the predefined initial state  $|\psi_3^{\text{in}}\rangle = |000\rangle$ , do not produce the expected target state  $|\psi_3\rangle$ , but instead produce a different state:  $|\psi_3'\rangle = \frac{1}{2}(|010\rangle + |011\rangle + \sqrt{2}|100\rangle)$ . Nevertheless, since this state is represented and labeled by the targeted SWEET state  $|\psi_3\rangle$ ,

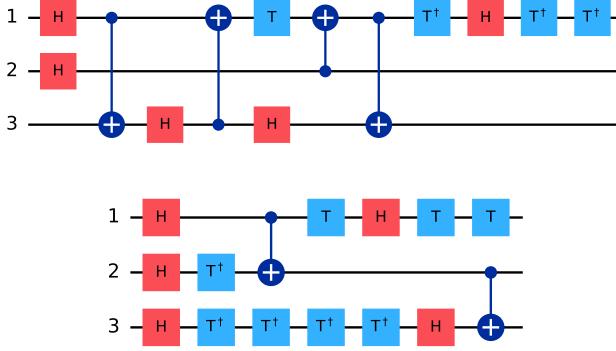


Fig. 4. Optimized quantum circuits for preparing three-qubit states. The figure compares two circuit designs, each optimized according to a different criterion. Top Panel: Displays the circuit produced when the optimization prioritizes a minimal gate count. This configuration results in a circuit with 13 gates and a depth of 11. Bottom Panel: Shows the outcome when a depth penalty is included in the optimization, directing the agent to minimize circuit depth. This strategy yields a circuit with a depth of 7 at the expense of a modestly increased gate count (15 gates). These results highlight the effectiveness of dynamic penalty strategies in circuit synthesis, enabling flexible prioritization between gate count and depth according to user-defined objectives. For further methodological details, refer to Sec.IV-K.

the algorithm is still considered successful. A more detailed discussion is provided in Appendix A.

## V. CONCLUSION

The efficient synthesis of quantum circuits remains a critical challenge to fully harness the potential of quantum computing. This work introduces a hybrid reward-driven reinforcement learning framework that successfully navigates the vast and sparse search space of quantum operations to discover resource-efficient circuits.

The reinforcement learning method involves a hybrid reward system that synergistically combines a precomputed static reward with dynamic penalties generated during the training phase. The static component creates a "breadcrumb trail" toward the target state due to its layered structure. Concurrently, dynamic penalties provide the flexibility to introduce user-defined constraints, such as minimizing circuit depth or costly T-gate counts, steering the agent toward genuinely efficient solutions.

To avoid wandering through the exponentially large space of state-action pairs, the Q-learning framework includes both exploration and exploitation of the space. To effectively balance between them, we used the well-known  $\epsilon$ -Greedy method, which efficiently guided the agent to the target regardless of the enormous space dimension.

To render the problem tractable and usefully discretize the state space, we focused our framework on States With Equal-amplitude and Encoded-phase Terms (SWEET states). We benchmarked our algorithm by producing circuits for graph states, minimizing circuit depth. The circuits constructed by the algorithm comprise the theoretically predicted minimum depth, according to the Vizing Theorem [40]. This success

demonstrates the algorithm ability to learn and exploit parallelizable gate structures implicitly, from the reward signal alone.

The versatility of the framework is further confirmed by synthesizing SWEET states using a universal gate set, proving its adaptability beyond structured problems. The performance of the algorithm is successfully tested for constructing circuits acting on 7 and 3-qubit states when the considered number of phase qubits is 1 and 3, respectively. The sparse matrix storage and layered reward shaping make the approach scalable to larger systems while keeping computational overhead manageable. Although the examples presented here involve minimizing the depth of the circuit, the very same penalty-based method can be used to penalize circuits with other faults, such as reducing the number of magic gates or entangling gates. Importantly, this work highlights how reinforcement learning can integrate classical reward shaping and quantum circuit constraints to navigate the challenges posed by the vastness of state space.

The proposed algorithm lends itself to a wide range of applications, both within the context of quantum circuit optimization and beyond. One particularly promising direction involves converting it into a circuit synthesizer, targeting unitary transformations rather than a fixed target state. This would leverage the fact that a unitary operator is fully characterized by its action on the computational basis, opening possibilities for discovering minimum-depth decompositions of known quantum gates.

A more immediate extension of our current algorithm concerns the analysis of the penalty patterns that we produce during the training part. Specifically, the  $R_{\text{dyn}}$  values constitute penalty maps, which are generated for a given target state. These maps could be generalized to construct informative initial penalty patterns when dealing with larger systems. This approach aims to provide a more scalable variant of the technique by establishing a structured initialization of the agent policy, potentially improving performance for higher qubit counts by accelerating the convergence of the learning procedure.

## ACKNOWLEDGMENT

The authors acknowledge support from Spanish MICIN grant PID2021-122547NB-I00 and the "MADQuantum-CM" project funded by Comunidad de Madrid (Programa de acciones complementarias) and by the Ministry for Digital Transformation and of Civil Service 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 Next Generation EU within the framework of the Digital Spain 2026 Agenda, the CAM Programa TEC-2024/COM-84 QUITEMAD-CM. This work has been financially supported by the project MADQuantum-CM, funded by Comunidad de Madrid (Programa de Acciones Complementarias) and by the Recovery, Transformation and Resilience Plan—Funded by the European Union—(NextGeneration

EU, PRTR-C17.II). KS and SG acknowledges support from the project MADQuantum-CM, funded by Comunidad de Madrid (Programa de Acciones Complementarias) and by the Recovery, transformation and Resilience Plan—Funded by the European Union—(NextGeneration EU, PRTR-C17.II). M.A. M.-D. has been partially supported by the U.S. Army Research Office Through Grant No. W911NF-14-1-0103.

## APPENDIX A TARGETING SWEET STATES

As we mentioned in Subsections II-A and IV-L, due to the infeasibility of labeling and representing each state of a given Hilbert space in an efficient computational way, we have restricted ourselves to the set of SWEET states. As a result, our algorithm can only register states having the SWEET states form. However, due to the inclusion of non-trivial gates in the algorithm, an action over a SWEET state can bring to a non-SWEET one. For instance, if we apply the Hadamard gate on the first qubit of the SWEET state,  $|\psi\rangle = \frac{1}{\sqrt{3}}(|000\rangle + |001\rangle + |010\rangle)$ , the final state we get is  $|\psi'\rangle = \frac{1}{\sqrt{6}}(2|000\rangle + |010\rangle + |011\rangle)$ , which is not SWEET anymore.

If we want to build the algorithm for such nontrivial gates, still being restricted to the set of SWEET states,  $\sum_j \alpha_j |x_j\rangle$  (see Eq.(1)), we must consider each SWEET state to represent a class of states,  $\sum_j a_j \alpha_j |x_j\rangle$ , where  $a_j$  are positive real numbers for all  $j$ s. In these cases, the algorithm treats any state falling outside the SWEET set but within a SWEET class as equivalent to the representative SWEET state of that class. As a result, in such situations, the algorithm does not target a specific state, but rather aims at the class of states represented by a given SWEET state, using the representative SWEET state as the target state, and constructs a circuit that produces a state belonging to that class.

Therefore, since Hadamard gates were used in Sec. IV-I, the algorithm, although targeting the SWEET state  $|\psi_3\rangle$ , constructed a circuit that does not produce that exact state. Instead, it generates a state that is represented by  $|\psi_3\rangle$  and belongs to the class labeled by it.

## REFERENCES

- [1] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 3rd ed. USA: Prentice Hall Press, 2009.
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: A Bradford Book, 2018.
- [3] M. Wiering and M. van Otterlo, *Reinforcement Learning: State-of-the-Art*. Springer Publishing Company, Incorporated, 2014.
- [4] S. Giordano and M. A. Martin-Delgado, “Reinforcement-learning generation of four-qubit entangled states,” *Phys. Rev. Res.*, vol. 4, p. 043056, Oct 2022.
- [5] I. Moflic and A. Paler, “Towards faster reinforcement learning of quantum circuit optimisation: Exponential reward functions,” in *Proceedings of the 18th ACM International Symposium on Nanoscale Architectures*, ser. NANOARCH ’23. New York, NY, USA: Association for Computing Machinery, 2024.
- [6] C. J. C. H. Watkins and P. Dayan, “Q-learning,” *Machine Learning*, vol. 8, no. 3, pp. 279–292, 1992.
- [7] A. Y. Ng, D. Harada, and S. J. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *Proceedings of the Sixteenth International Conference on Machine Learning*, ser. ICML ’99. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1999, p. 278–287.
- [8] H. S. Jakab and L. Csató, “Sparse approximations to value functions in reinforcement learning,” in *Artificial Neural Networks*, P. Koprinkova-Hristova, V. Mladenov, and N. K. Kasabov, Eds. Cham: Springer International Publishing, 2015, pp. 295–314.
- [9] G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester, “Challenges of real-world reinforcement learning: definitions, benchmarks and analysis,” *Mach. Learn.*, vol. 110, no. 9, p. 2419–2468, Sep. 2021.
- [10] D. Schlingemann and R. F. Werner, “Quantum error-correcting codes associated with graphs,” *Phys. Rev. A*, vol. 65, p. 012308, Dec 2001.
- [11] R. Raussendorf, D. E. Browne, and H. J. Briegel, “Measurement-based quantum computation on cluster states,” *Phys. Rev. A*, vol. 68, p. 022312, Aug 2003.
- [12] M. Kafatos, Ed. *Bell’s Theorem, Quantum Theory and Conceptions of the Universe*, ser. Fundamental Theories of Physics. Dordrecht: Springer Dordrecht, 1989, vol. 37.
- [13] W. Dür, G. Vidal, and J. I. Cirac, “Three qubits can be entangled in two inequivalent ways,” *Phys. Rev. A*, vol. 62, p. 062314, Nov 2000.
- [14] D. Gottesman, “The heisenberg representation of quantum computers,” in *Group22: Proceedings of the XXII International Colloquium on Group Theoretical Methods in Physics*, S. P. Corney, R. Delbourgo, and P. D. Jarvis, Eds. Cambridge, MA: International Press, 1999, pp. 32–43.
- [15] M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, 2010.
- [16] A. Galindo and M. A. Martín-Delgado, “Information and computation: Classical and quantum aspects,” *Rev. Mod. Phys.*, vol. 74, pp. 347–423, May 2002.
- [17] F. G. Brandão, W. Chemissany, N. Hunter-Jones, R. Kueng, and J. Preskill, “Models of quantum complexity growth,” *PRX Quantum*, vol. 2, p. 030316, Jul 2021.
- [18] J. Watrous, *Quantum Computational Complexity*. New York, NY: Springer New York, 2009, pp. 7174–7201.
- [19] B. Craps, M. D. Clerck, O. Evnin, and P. Hacker, “Integrability and complexity in quantum spin chains,” *SciPost Phys.*, vol. 16, p. 041, 2024.
- [20] S. S. Haque, C. Jana, and B. Underwood, “Saturation of thermal complexity of purification,” *Journal of High Energy Physics*, vol. 2022, no. 1, p. 159, 2022.
- [21] P. Caputa and S. Liu, “Quantum complexity and topological phases of matter,” *Phys. Rev. B*, vol. 106, p. 195125, Nov 2022.
- [22] M. A. Nielsen, “A geometric approach to quantum circuit lower bounds,” *Quantum Info. Comput.*, vol. 6, no. 3, p. 213–262, May 2006.
- [23] M. A. Nielsen, M. R. Dowling, M. Gu, and A. C. Doherty, “Quantum computation as geometry,” *Science*, vol. 311, no. 5764, pp. 1133–1135, 2006.
- [24] R. A. Jefferson and R. C. Myers, “Circuit complexity in quantum field theory,” *Journal of High Energy Physics*, vol. 2017, no. 10, p. 107, 2017.
- [25] L. Hackl and R. C. Myers, “Circuit complexity for free fermions,” *Journal of High Energy Physics*, vol. 2018, no. 7, p. 139, 2018.
- [26] T. Ali, A. Bhattacharyya, S. S. Haque, E. H. Kim, and N. Moynihan, “Time evolution of complexity: a critique of three methods,” *Journal of High Energy Physics*, vol. 2019, no. 4, p. 87, 2019.
- [27] R. Elmasri and S. Navathe, *Fundamentals of Database Systems*, 6th ed. USA: Addison-Wesley Publishing Company, 2010.
- [28] G. A. Rummery and M. Niranjan, “On-line q-learning using connectionist systems,” in *Technical Report CUED/F-INFENG/TR 166, Cambridge University Engineering Department*, 1994.
- [29] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley, “A survey of multi-objective sequential decision-making,” *Journal of Artificial Intelligence Research*, vol. 48, no. 1, p. 67–113, Oct. 2013.
- [30] S. Elfwing and B. Seymour, “Parallel reward and punishment control in humans and robots: Safe reinforcement learning using the maxpain algorithm,” in *2017 Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, 2017, pp. 140–147.
- [31] S. Devlin and D. Kudenko, “Theoretical considerations of potential-based reward shaping for multi-agent systems,” in *The 10th International Conference on Autonomous Agents and Multiagent Systems - Volume*

- I*, ser. AAMAS '11. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2011, p. 225–232.
- [32] R. S. Sutton, “Dyna, an integrated architecture for learning, planning, and reacting,” *SIGART Bull.*, vol. 2, no. 4, p. 160–163, Jul. 1991.
- [33] ———, “Integrated architectures for learning, planning, and reacting based on approximating dynamic programming,” in *Machine Learning Proceedings 1990*, B. Porter and R. Mooney, Eds. San Francisco (CA): Morgan Kaufmann, 1990, pp. 216–224.
- [34] C. Blundell, B. Uria, A. Pritzel, Y. Li, A. Ruderman, J. Z. Leibo, J. Rae, D. Wierstra, and D. Hassabis, “Model-free episodic control,” 2016.
- [35] V. Mnih, K. Kavukcuoglu, D. Silver *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [36] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Advances in Neural Information Processing Systems*, ser. NIPS'99, S. Solla, T. Leen, and K. Müller, Eds., vol. 12. Cambridge, MA, USA: MIT Press, 1999, p. 1057–1063.
- [37] V. Konda and J. Tsitsiklis, “Actor-critic algorithms,” in *Advances in Neural Information Processing Systems*, S. Solla, T. Leen, and K. Müller, Eds., vol. 12. MIT Press, 1999.
- [38] M. Hein, W. Dür, J. Eisert, R. Raussendorf, M. Van den Nest, and H.-J. Briegel, “Entanglement in graph states and its applications,” in *Proceedings of the International School of Physics “Enrico Fermi”: Quantum Computers, Algorithms and Chaos*, ser. Enrico Fermi School of Physics, G. Casati and S. Montangero, Eds. Varenna, Italy: IOS Press, 2006.
- [39] M. Hein, J. Eisert, and H. J. Briegel, “Multiparty entanglement in graph states,” *Phys. Rev. A*, vol. 69, p. 062311, Jun 2004.
- [40] J. Misra and D. Gries, “A constructive proof of vizing’s theorem,” *Information Processing Letters*, vol. 41, no. 3, pp. 131–133, 1992.
- [41] C. Szepesvari, *Algorithms for Reinforcement Learning*. Morgan and Claypool Publishers, 2010.
- [42] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [43] M. Amy, D. Maslov, M. Mosca, and M. Roetteler, “A meet-in-the-middle algorithm for fast synthesis of depth-optimal quantum circuits,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 32, no. 6, pp. 818–830, 2013.
- [44] O. Golubitsky and D. Maslov, “A study of optimal 4-bit reversible toffoli circuits and their synthesis,” *IEEE Transactions on Computers*, vol. 61, no. 9, pp. 1341–1353, 2012.