

Simultaneous Discovery of Quantum Error Correction Codes and Encoders with a Noise-Aware Reinforcement Learning Agent

Jan Olle,^{1,*} Remmy Zen,¹ Matteo Puviani,¹ and Florian Marquardt^{1,2}

¹*Max Planck Institute for the Science of Light, Staudtstraße 2, 91058 Erlangen, Germany*

²*Department of Physics, Friedrich-Alexander Universität Erlangen-Nürnberg, Staudtstraße 5, 91058 Erlangen, Germany*

(Dated: November 9, 2023)

Finding optimal ways to protect quantum states from noise remains an outstanding challenge across all quantum technologies, and quantum error correction (QEC) is the most promising strategy to address this issue. Constructing QEC codes is a complex task that has historically been powered by human creativity with the discovery of a large zoo of families of codes. However, in the context of real-world scenarios there are two challenges: these codes have typically been categorized only for their performance under an idealized noise model and the implementation-specific optimal encoding circuit is not known. In this work, we train a Deep Reinforcement Learning agent that automatically discovers both QEC codes and their encoding circuits for a given gate set, qubit connectivity, and error model. We introduce the concept of a noise-aware meta-agent, which learns to produce encoding strategies simultaneously for a range of noise models, thus leveraging transfer of insights between different situations. Moreover, thanks to the use of the stabilizer formalism and a vectorized Clifford simulator, our RL implementation is extremely efficient, allowing us to produce many codes and their encoders from scratch within seconds, with code distances varying from 3 to 5 and with up to 20 physical qubits. Our approach opens the door towards hardware-adapted accelerated discovery of QEC approaches across the full spectrum of quantum hardware platforms of interest.

I. INTRODUCTION

There is an ongoing global effort to develop a new generation of quantum technologies with an unprecedented level of control over individual quantum states of many-particle quantum systems. This field encompasses four areas [1]: quantum communication, simulation, computation and sensing, each of them promising to drastically improve on preceding classical technologies. An outstanding challenge that is present in all the aforementioned areas is that quantum states are vulnerable to the unwanted effects of noise; if not addressed, the advantages offered over classical technologies disappear altogether.

Quantum error correction (QEC) is a field which emerges from the union of quantum mechanics and classical error correction [2], and it is the approach that is thought to be essential to achieve maturity in the current wave of quantum technologies. The core idea of QEC is to redundantly embed the quantum information within a subspace (called *code space*) of a *larger* Hilbert space in such a way that different errors map the code space to mutually orthogonal subspaces. If successful, the action of each of these errors can be reverted and the quantum process can continue error-free [3]. QEC is *necessary* for large-scale fault-tolerant quantum computing [4, 5]. The past few years have witnessed dramatic progress in experimental realizations of QEC on different platforms [6–10], reaching a point where the lifetime of qubits has been extended by applying QEC [11].

Since Shor’s original breakthrough [12], different qubit-based QEC codes have been constructed, both analyti-

cally and numerically, leading to a zoo of codes, each of them conventionally labeled $[[n, k, d]]$, where n is the number of physical qubits, k the number of encoded logical qubits, and d the code distance that defines the number $d - 1$ of detectable errors. The first examples are provided by the $[[5, 1, 3]]$ perfect code [13], the $[[7, 1, 3]]$ Steane [14] and the $[[9, 1, 3]]$ Shor [12] codes, which encode one logical qubit into 5, 7 and 9 physical qubits, respectively, being able to detect up to 2 physical errors and correct up to 1 error on any physical qubit. The most promising approach so far is probably the family of the so-called *toric* or *surface codes* [15], which encode a logical qubit into the joint entangled state of a $d \times d$ square of physical qubits. More recently, examples of quantum Low-Density Parity Check (LDPC) codes that are competitive with the surface code have been discovered [16].

Numerical techniques have already been employed to construct QEC codes. Often, this has involved greedy algorithms, which may lead to sub-optimal solutions but can be relatively fast. For instance, in [17] a greedy algorithm was implemented to extend codes to higher distance, code concatenation was explored in [18], and greedy search for finding stabilizer codes was used in [19]. Often, such numerical methods for QEC code construction were restricted to finding a subclass of codes of a particular structure, e.g. using reduction to classical code search [20].

With the recent advent of powerful tools from the domains of machine learning and, more generally, Artificial Intelligence (AI), there are new opportunities for the automated discovery of QEC schemes. Of special relevance to this work is the realm of Reinforcement Learning (RL), which is designed to solve complex decision-making problems by autonomously following an action-reward scheme [21]. The task to solve is encoded in a

* jan.olle@mpl.mpg.de

reward function, and the aim of RL training algorithms is to maximize such a reward over time. RL can provide new answers to difficult questions, in particular in fields where optimization in a high-dimensional search space plays a crucial role. For this reason, reinforcement learning can be an efficient tool to deal with the task of QEC code construction and encoding.

The first example of RL-based automated discovery of QEC strategies [22] did not rely on any human knowledge of QEC concepts. While this allowed exploration without any restrictions, e.g. going beyond stabilizer codes, it was limited to only small qubit numbers. More recent works have moved towards optimizing only certain QEC subtasks, injecting substantial human knowledge. For example, RL has been used for optimization of given QEC codes [23], and to discover tensor network codes [24] or codes based on "Quantum Lego" parametrizations [25, 26]. Additionally, RL has been used to find efficient decoding processes [27–29].

At this moment, a multitude of experimental platforms are scaling up towards the regime of qubit numbers that make it possible to implement QEC (this includes especially various superconducting qubit architectures, ion traps, quantum dots, and neutral atoms). Given the strong differences in native gate sets, qubit connectivities, and relevant noise models, there is a strong need for a flexible and efficient scheme to automatically discover not only codes but also efficient encoding circuits, adapted to the platform at hand. In our work, we implement a scheme based on deep RL in order to simultaneously discover QEC codes together with the encoding circuit from scratch, tailored to specific noise models, native gate sets and connectivities, minimizing the circuit size for improved hardware efficiency. In particular, our RL agent can be made *noise-aware*, meaning that one and the same agent is able to switch its encoding strategy based on the specific noise that is present in the system. Our approach is flexible and general, as well as efficient through the use of a parallelized Clifford circuit simulator.

While [30] also set themselves the task of finding both codes and their encoding circuits, this was done using variational quantum circuits involving continuously parametrized gates, which leads to much more costly numerical simulations and eventually only an approximate QEC scheme. By contrast, our RL-based approach does not rely on any human-provided circuit ansatz, can use directly any given discrete gate set, is able to exploit highly efficient Clifford simulations, and produces a meta-agent able to cover strategies for a range of noise models.

The paper is organized as follows: in Section II we provide the theoretical background for stabilizer codes, code classification and reinforcement learning. In Section III we describe our approach to build a noise-aware reinforcement learning agent that discovers multiple QEC codes in asymmetric noise channels. Finally, we present and analyze our numerical results in Section IV.

II. BACKGROUND

A. Stabilizer Codes in Symmetric and Asymmetric Error Channels

Some of the most promising QEC codes are based on the stabilizer formalism [31], which leverages the properties of the Pauli group G_n on n qubits. The basic idea of the stabilizer formalism is that many quantum states of interest for QEC can be more compactly described by listing the set of n operators that *stabilize* them, where an operator O stabilizes a state $|\psi\rangle$ if $O|\psi\rangle = |\psi\rangle$. The Pauli group on a single qubit G_1 is defined as the set of Pauli matrices X, Y, Z with the overall phases $\pm 1, \pm i$, which form a group under matrix multiplication. The generalization to n qubits consists of all n -fold tensor products of Pauli matrices (called *Pauli strings*). For the purposes of QEC, global phases can be ignored.

A code that encodes k logical qubits into n physical qubits is a 2^k -dimensional subspace (the *code space* \mathcal{C}) of the full 2^n -dimensional Hilbert space. It is completely specified by the set of Pauli strings $S_{\mathcal{C}}$ that *stabilize* it, i.e. $S_{\mathcal{C}} = \{s_i \in G_n \mid s_i|\psi\rangle = |\psi\rangle, \forall |\psi\rangle \in \mathcal{C}\}$. $S_{\mathcal{C}}$ is called the stabilizer group of \mathcal{C} and is usually written in terms of its group generators g_i as $S_{\mathcal{C}} = \langle g_1, g_2, \dots, g_{n-k} \rangle$.

Noise affecting quantum processes can be represented using the so-called *operator-sum* representation [32], where a quantum noise channel \mathcal{N} induces dynamics on the state ρ according to

$$\mathcal{N}(\rho) = \sum_{\alpha} E_{\alpha} \rho E_{\alpha}^{\dagger}, \quad (1)$$

where E_{α} are *Kraus operators*, satisfying $\sum_{\alpha} E_{\alpha}^{\dagger} E_{\alpha} = I$. The most elementary example is the so-called *depolarizing* noise channel,

$$\mathcal{N}_{\text{DP}}(\rho) = p_I \rho + p_X X \rho X + p_Y Y \rho Y + p_Z Z \rho Z, \quad (2)$$

where $p_I + p_X + p_Y + p_Z = 1$ and the set of Kraus operators are $E_{\alpha} = \{\sqrt{p_I} I, \sqrt{p_X} X, \sqrt{p_Y} Y, \sqrt{p_Z} Z\}$. When considering n qubits, one can generalize the depolarizing noise channel by introducing the *global* depolarizing channel,

$$\mathcal{N}_{\text{GDP}}(\rho) = \bigotimes_{j=1}^n \mathcal{N}_{\text{DP}}^{(j)}(\rho_j), \quad (3)$$

consisting of local depolarizing channels acting on each qubit j independently. Taken as is, this error model generates *all* 4^n Pauli strings by expanding (3). A commonly used simplification is the following. Assume that all error probabilities are identical, i.e. $p_X = p_Y = p_Z \equiv p$ (and $p_I = 1 - 3p$). Then, the probability that a given error occurs decreases with the number of qubits it affects. For instance, if we consider 3 qubits, the probability associated with XII is $p(XII) = p(1 - 3p)^2$, and in general the leading order contribution to the probability of an error affecting m qubits is p^m . This leads to the concept of the

weight of an operator as the number of qubits on which it differs from the identity and to a hierarchical approach to building QEC codes. In particular, stabilizer codes are described by specifying what is the minimal weight in the Pauli group that they cannot detect.

The fundamental theorem in QEC is a set of necessary and sufficient conditions for quantum error detection discovered independently by Bennett, DiVincenzo, Smolin and Wootters [33], and by Knill and Laflamme in [34] (KL conditions from now on). These state that a code \mathcal{C} with associated stabilizer group $S_{\mathcal{C}}$ can *detect* a set of errors $\{E_{\mu}\} \subseteq G_n$, if and only if for all E_{μ} we have either

$$\{E_{\mu}, g_i\} = 0, \quad (4)$$

for at least one g_i , or the error itself is harmless, i.e.

$$E_{\mu} \in S_{\mathcal{C}}. \quad (5)$$

The smallest weight in G_n for which none of the above two conditions hold is called the *distance* of the code. For instance, a distance-3 code is capable of detecting *all* Pauli strings of up to weight 2, meaning that KL conditions (4), (5) are satisfied for *all* Pauli strings of weights 0, 1 and 2. Moreover, the *smallest* weight for which these are not satisfied is 3, meaning that there is *at least one* weight-3 Pauli string violating both (4) and (5). However, *some* weight-3 Pauli strings (and higher weights) will satisfy the KL conditions, in general.

While these conditions are framed in the context of quantum error detection, there is a direct correspondence with quantum error correction. Indeed, a quantum code of distance d can *correct* all errors of up to weight $t = \lfloor (d-1)/2 \rfloor$ [31]. Thus, $d = 3$ codes can both detect up to *all* weight-2 errors and can correct *all* weight-1 errors. If all the errors that are detected with a weight smaller than d obey (4), the code is called *non-degenerate*. On the other hand, if some of the errors satisfy (5), the code is called *degenerate*.

The default weight-based $[[n, k, d]]$ classification of QEC codes implicitly assumes that the error channel is symmetric, meaning that the probabilities of Pauli X, Y and Z errors are equal. However, this is usually not the case in experimental setups: for example, dephasing (Z errors) may dominate bit-flip errors.

In our work, we will consider an asymmetric noise channel where $p_X = p_Y$ but $p_X \neq p_Z$. To quantify the asymmetry, we use the bias parameter c_Z [30], defined as

$$c_Z = \frac{\log p_Z}{\log p_X}. \quad (6)$$

For symmetric error channels, $c_Z = 1$. If Z-errors dominate, then $0 < c_Z < 1$, since $p_Z = p_X^{c_Z}$ and $p_X, p_Z \ll 1$; conversely $c_Z > 1$ when X/Y errors are more likely than Z errors.

The weight of operators and the code distance can both be generalized to asymmetric noise channels [35–38]. Consider a Pauli string operator E_{μ} and denote as w_X the number of Pauli X inside E_{μ} (likewise for Y, Z).

Then one can introduce the c_Z -effective weight [30] of E_{μ} as

$$w_e(E_{\mu}, c_Z) = w_X(E_{\mu}) + w_Y(E_{\mu}) + c_Z w_Z(E_{\mu}), \quad (7)$$

which reduces to the symmetric weight for $c_Z = 1$, as expected. The c_Z -effective distance of a code $d_e(c_Z)$ is then defined [30] as the largest possible integer such that the KL conditions (4), (5) hold for all Pauli strings E_{μ} with $w_e(E_{\mu}, c_Z) < d_e(c_Z)$. Like in the symmetric noise case, the meaning of this effective distance is that *all* error operators with an effective weight smaller than d_e can be detected.

B. Code Classification

It is well known that there is no unique way to describe quantum codes. For instance, there are multiple sets of code generators that generate the same stabilizer group, hence describing the *same* code. Moreover, the choice of logical basis is not unique and qubit labeling is arbitrary. While such redundancies are convenient for describing quantum codes in a compact way, comparing and classifying different codes can be rather subtle. Fortunately, precise notions of code equivalence have been available in the literature since the early days of this field. In this work, we will refer to *families* of codes based on their quantum weight enumerators (QWE) [39], $A(z)$ and $B(z)$, which are polynomials with coefficients

$$\begin{aligned} A_j &= \frac{1}{(2^k)^2} \sum_{w(E_{\mu})=j} \text{Tr}(E_{\mu} P_{\mathcal{C}}) \text{Tr}(E_{\mu}^{\dagger} P_{\mathcal{C}}), \\ B_j &= \frac{1}{2^k} \sum_{w(E_{\mu})=j} \text{Tr}(E_{\mu} P_{\mathcal{C}} E_{\mu}^{\dagger} P_{\mathcal{C}}), \end{aligned} \quad (8)$$

where j runs from 0 to n and $P_{\mathcal{C}}$ is the orthogonal projector onto the code space. Intuitively, A_j counts the number of error operators of weight j in $S_{\mathcal{C}}$ while B_j counts the number of error operators of weight j that commute with all elements of $S_{\mathcal{C}}$. Logical errors are thus the ones that commute with $S_{\mathcal{C}}$ but are not in $S_{\mathcal{C}}$, and these are counted with $B_j - A_j$.

Such a classification is especially useful in scenarios with symmetric noise channels, where it is irrelevant whether the undetected errors contain a specific Pauli operator at a specific position. However, such a distinction can in principle be important in asymmetric noise channels. While generalizing (8) to asymmetric noise channels with effective weights is straightforward, comparing codes across different values of noise bias becomes cumbersome. Hence, in the present work we will always refer to (symmetric) code families according to (8) for all values of c_Z , i.e. we will effectively pretend that $c_Z = 1$ when computing the weight enumerators of asymmetric codes.

C. Reinforcement Learning

Reinforcement Learning (RL) [40] is designed to discover optimal action sequences in decision-making problems. The goal in any RL task is encoded by choosing a suitable *reward* r , a quantity that measures how well the task has been solved, and consists of an *agent* (the entity making the decisions) interacting with an *environment* (the physical system). In each time step t , the environment's state s_t is observed. Based on this observation, the agent takes an action a_t which then affects the current state of the environment. A *trajectory* is a sequence of state and action pairs that the agent takes. An *episode* is a trajectory from an initial state to a terminal state. For each action, the agent receives a reward r_t , and the goal of RL algorithms is to maximize the expected cumulative reward (return), $\mathbb{E}[\sum_t r_t]$. The agent's behavior is defined by the *policy* $\pi_\theta(a_t|s_t)$, which denotes the probability of choosing action a_t given observation s_t , and is parameterized by a neural network with parameters θ .

Within RL, policy gradient methods [21] optimize the policy by maximizing the expected return with respect to the parameters θ with gradient ascent. One of the most successful algorithms within policy gradient methods is the actor-critic algorithm [41]. The idea is to have two neural networks: an actor network that acts as the agent and that defines the policy, and a critic network, which measures how good was the action taken by the agent. In this paper, we use a state-of-the-art policy-gradient actor-critic method called Proximal Policy Optimization (PPO) [42], which improves the efficiency and stability of policy gradient methods.

III. REINFORCEMENT LEARNING APPROACH TO QEC CODE DISCOVERY

The main objective of this work is to automatize the discovery of QEC codes and their encoding circuits using RL. We will be interested in a scenario where the encoding circuit is assumed to be error-free. This is applicable to quantum communication or quantum memories, where the majority of errors happen during transmission over a noisy channel or during the time the memory is retaining the information.

Eventually we will show that it is possible to train a meta-agent that is capable of adapting its strategy according to the noise model, without any retraining. This leverages the concept of transfer learning, where improvements gained in training for one scenario (here, one value of a noise parameter) carry over and accelerate the training progress for other scenarios. A scheme of our approach can be found in Fig. 1.

A. Encoding Circuit

In order to encode the state of k logical qubits on n physical qubits one must find a sequence of quantum gates that will entangle the quantum information in such a way that QEC is possible with respect to a target noise channel. Initially, we imagine the first k qubits as the original containers of our (yet unencoded) quantum information, which can be in any state $|\psi\rangle \in (\mathbb{C}_2)^{\otimes k}$. The remaining $n - k$ qubits are chosen to each be initialized in the state $|0\rangle$. These will be turned into the corresponding logical state $|\psi\rangle_L \in (\mathbb{C}_2)^{\otimes n}$ via the application of a sequence of Clifford gates on any of the n qubits (where Clifford gates are defined to be those that map Pauli strings to Pauli strings and are generated by the Hadamard H , the Phase S and the CNOT gates). In the stabilizer formalism, this means that initially the generators of the stabilizer group are

$$Z_{k+1}, Z_{k+2}, \dots, Z_n. \quad (9)$$

The task of the RL agent is to discover a suitable encoding sequence of gates. After applying each gate, the $n - k$ code generators (9) are updated. The agent then receives a representation of these generators as input (as its observation) and suggests the next gate (action) to apply. In this way, an encoding circuit is built up step by step, taking into account the available gate set and connectivity for the particular hardware platform. This process terminates when the KL conditions (4), (5) are satisfied for the target error channel and the learned circuit can then be used to encode *any* state $|\psi\rangle$ of choice.

Before discussing how to implement the KL conditions in terms of a reward function, we illustrate this process with the encoding of the three-qubit repetition code, which is a $[[3, 1]]$ code that can correct single-qubit bit flips. This means that there are two code generators and that the targeted error set consists of all weight-1 and weight-2 bit-flip (Pauli- X) operators ($XII, IXI, IIX, XXI, XIX, IXX$), since detecting this set implies being able to correct all single-qubit bit-flips. The generators start being IZI, IIZ in the absence of gates, and transform into ZZI, ZIZ after applying two CNOT gates with control on the first qubit and target on the second and third qubits. Initially, only the error operator XII is not detected. After the first CNOT gate, the undetected error becomes XXI , and after the second CNOT all target errors are detected.

B. Reward

The most delicate matter in RL problems is building a suitable reward for the task at hand. Our goal is to design an agent that, given an error model that includes a set of errors $\{E_\mu\} \subseteq G_n$ with associated occurrence probabilities $\{p_\mu\}$, is able to find an encoding sequence that protects the quantum information from such noise.

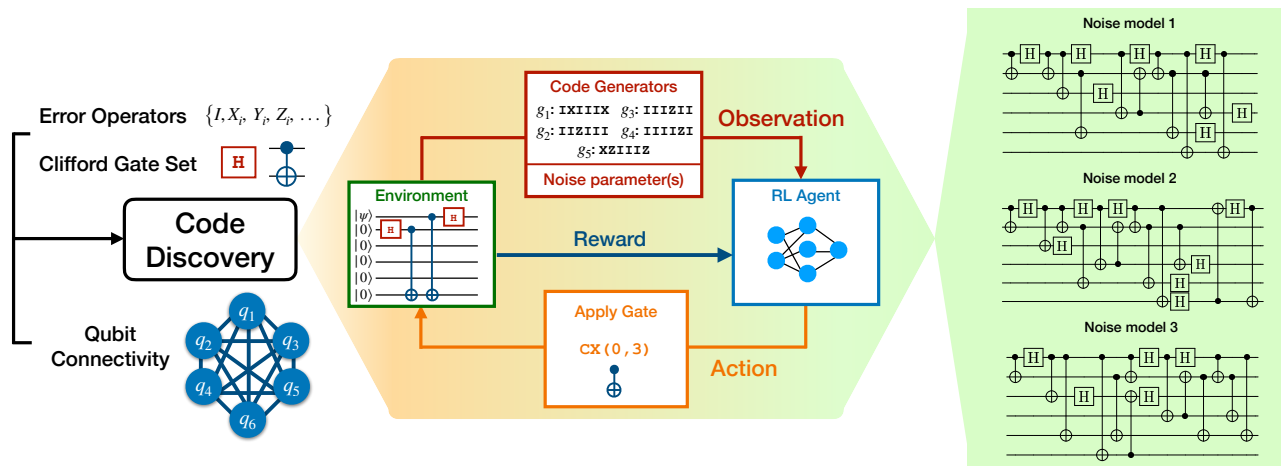


FIG. 1. QEC code and encoding discovery using a noise-aware RL meta-agent. A set of error operators, a gate set and qubit connectivity is chosen. Different error models can be considered by varying some noise parameters, which are fed as an observation to the agent. The agent then builds a circuit using the available gate set and connectivity that detects the most likely errors from the target error model by using a reward based on the Knill-Laflamme QEC conditions according to Eq. (10). After training, a single RL agent is able to find suitable encodings for different noise models, which are able to encode *any* state $|\psi\rangle$ of choice.

Let us consider a quantum communication setup to have a concrete picture. Here, we imagine Alice and Bob exchanging some quantum bits of information contained in a state $|\psi\rangle$. We assume that they are able to implement all gates and measurements without errors and that errors only happen while the message is traveling through the communication channel. Alice encodes state $|\psi\rangle$ into $|\psi\rangle_L$ and sends it through the noisy channel, after which it is received by Bob in the form of a possibly corrupted, mixed state. Since Bob knows the encoding that Alice has used, he also knows what are the stabilizer generators. Hence, Bob proceeds to perform syndrome measurements according to the stabilizer generators of the code and finally corrects the errors that he believes have happened along the way. If done perfectly, the corrected state becomes the original state $|\psi\rangle_L$ sent by Alice. However, even when both Alice and Bob have access to perfect gates and measurements, the probability of recovering the state is not one. The reason is that multiple errors may trigger the same syndrome measurement, and in some cases Bob will mistakenly correct for an error that did not actually happen. An option for our RL agent could thus be to maximize the probability of recovering the initial state, or what is the same, minimizing the probability that Bob applies a wrong error correction operation. Unfortunately, optimizing for this task is computationally very expensive. Indeed, for each syndrome, there are different errors that could have triggered it. On top of that, from all the errors that could have triggered that syndrome, computational resources have to be employed in finding the most dangerous one, which then has to be selected as the candidate error to be corrected. Such expensive calculations would have to be carried out in every step of the RL procedure.

A much cheaper alternative that avoids performing such error categorization is to use a scheme where the cumulative reward (which RL optimizes) simply is maximized whenever all the KL conditions are fulfilled. One implementation of this idea uses the weighted KL sum as an instantaneous reward:

$$r_t = - \sum_{\mu} \lambda_{\mu} K_{\mu}, \quad (10)$$

where $K_{\mu} = 0$ if either (4) or (5) are satisfied for the corresponding error operator E_{μ} , and $K_{\mu} = 1$ otherwise. Here λ_{μ} are real positive numbers that we will take to be hyperparameters. For any choice of λ_{μ} , if all errors in $\{E_{\mu}\}$ can be detected, the reward is zero, and is negative otherwise, thus leading the agent towards the goal we set out to solve. Later, we will also be interested in situations where not all errors can be corrected simultaneously and a good compromise has to be found. In that case, one simple heuristic choice for the reward (10) would be $\lambda_{\mu} = p_{\mu}$, giving more weight to errors that occur more frequently. While we will later see that maximizing the KL reward given here is not precisely equivalent to minimizing the overall probability of incorrectly classified errors, one can still expect a reasonable performance at this task, which would be the ultimate goal in any QEC scheme.

RL optimizes the cumulative reward R , i.e. the sum of r_t over all time steps. Therefore, defining the instantaneous reward r_t in terms of the weighted KL sum means we are asking the agent that it keeps this sum as small as possible *on average*. While this might seem less desirable than setting R itself to be the KL sum at the final time step, we have heuristically found that the ansatz adopted here produces good training behavior and is still able to

find codes and their encoding circuits.

Referring back to the encoding procedure for the 3-qubit repetition code, the weighted KL sum starts being $p_X p_I^2$ because only the error operator XII is not detected, changing to $p_X^2 p_I < p_X p_I^2$ after application of the first CNOT gate. When all KL conditions are satisfied, the weighted KL sum is zero and all error operators that were considered can be detected, leading to a successful encoding.

A nice feature of our reward is that one can favor certain types of codes. For instance, we can target non-degenerate codes only by ignoring (5). Moreover, making the reward non-positive favors short gate sequences, which is desirable when preparing these codes in actual quantum devices.

C. Noise-aware meta-agent

Regarding the error channel to be targeted, here there are in principle several choices that can be made. The most straightforward one is choosing a global depolarizing channel as given by (3). This still allows for asymmetric noise, i.e. different probabilities p_X, p_Y, p_Z . One option would be to train an agent for any given, fixed choice of these probabilities, necessitating retraining if these characteristics change. However, we want to go beyond that and discover a single agent being capable of deciding what is the optimal encoding strategy for *any* level of bias in the noise channel (6). For instance, we want this noise-aware agent to be able to understand that it should prioritize detecting ZZZ errors over XX errors in a scenario with $c_Z = 0.5$, but that it should do the opposite for larger values of c_Z . This translates into two aspects: The first one is that the agent has to receive the noise parameters as input. In the illustrative example further below, we will choose to supply the bias parameter c_Z as an extra observation, while keeping the overall error probability fixed. The second aspect is that the list of error operators will have to contain more operators than the total number that can actually be detected reliably, since it is now part of the agent’s task to prioritize some of those errors while ignoring the least likely errors. All in all, the list of operators participating in the reward (10) will be fixed and we will choose those with at most (symmetric) weight $d-1$, for a certain d ; the idea is then that when varying c_Z , the probabilities of these errors will change and the agent will have to figure out which are the most dangerous errors that must be detected in every case.

As we argued in Sec. II, if one expands (3), Pauli strings of all weights would participate in the error set $\{E_\mu\}$. Imposing a cutoff at (symmetric) weight $d-1$ implies that the number of error operators to keep track of is

$$|\{E_\mu\}|_{w \leq d-1} = \sum_{w=0}^{d-1} 3^w \binom{n}{w}, \quad (11)$$

which grows exponentially with d . As we will see, this exponential growth of the number of error operators that have to be tracked will impose the most severe limitation in our approach.

D. Fast parallelized Clifford simulator

RL algorithms exploit trial-and-error loops until a signal of a good strategy is picked up and convergence is reached, so it is of paramount importance that simulations of our RL environment are extremely fast. Luckily, thanks to the Gottesman-Knill theorem, the Clifford circuits needed here can be simulated efficiently classically. Optimized numerical implementations exist, e.g. STIM [43]. However, in an RL application we want to be able to run multiple agents in parallel in an efficient, vectorized way that is compatible with modern machine learning frameworks. For that reason, we have implemented our own special-purpose vectorized Clifford simulator. Briefly, we use the symplectic binary formalism of the Pauli group [44] to represent the stabilizer generators. S_C is then represented by a *check matrix* H [44], which is a $(n-k) \times 2n$ binary matrix where each row i represents the Pauli string g_i from S_C . Clifford gates are also implemented using binary matrices, and we achieve a massive parallelization by running Clifford-based simulations of quantum circuits in parallel, meaning that the agent interacts with a batch of RL environments (quantum circuits) at every timestep. Our Clifford simulator is implemented using JAX [45], a state-of-the-art modern machine learning framework with good vectorization and just-in-time compilation capabilities. On top of that, we also train multiple RL agents in parallel on a single GPU. This is achieved by interfacing with PUREJAXRL [46], a library that offers a high-performance end-to-end JAX RL implementation. We plan to make the codes of our implementation public in the near future.

The efficiency of Clifford simulations, as well as the use of powerful state-of-the-art RL algorithms, enables us to easily go beyond the recent results of [30], which were based on variational quantum circuits. In particular, we are able to straightforwardly discover codes and encoding circuits for both larger number of qubits (14 vs 20) and larger code distances (4 vs 5).

IV. RESULTS

We will first illustrate the basic workings of our approach for a symmetric noise channel before introducing the meta-agent that is able to simultaneously discover strategies for a range of noise models.

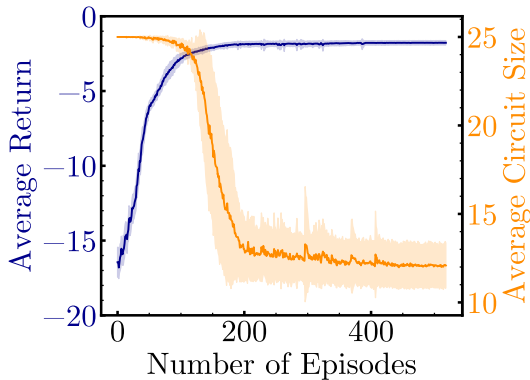


FIG. 2. Example of a training trajectory for $[[7, 1, 3]]$ code discovery. Here, 16 parallel agents each interact with batches of 32 circuits processed in parallel. Each agent finds a different encoding circuit, and the training finishes in 16 sec on a single GPU.

A. Codes in a symmetric depolarizing noise channel

We now illustrate the versatility of our approach by discovering a library of different codes and their associated encoding circuits. We fix the error model to be a symmetric depolarizing channel with error probability p , meaning $p_I = 1 - 3p$, $p_X = p_Y = p_Z = p$, and thus no noise parameter is needed. We also vary the target code distance d from 3 to 5. The corresponding target error set is $E_\mu = \{I, X_i, Y_i, Y_j, X_i X_j, \dots, Z_i Z_j\}$ for $d = 3$, and likewise for $d = 4, 5$, with the set for $d = 5$ including all Pauli string operators of up to weight 4. In our numerical experiments, we choose $p_I = 0.9$, which enters the reward function where we choose the weights proportional to the error probabilities (recall the discussion above). Although the final discovered codes and encoding circuits are independent of this choice, it can affect the learning progress.

For illustrative purposes, we take the gateset to be $\{H_i, \text{CNOT}(i < j)\}$, i.e. a *directed* all-to-all connectivity, which is sufficient given that our unencoded logical state is at the first qubits by design. Nevertheless, both using alternative gatesets or other connectivities works well. The final physical hyperparameter is the maximal number of gates that we allow before restarting the learning trajectory. This number will be varied from 20 to 50, depending on the target code parameters. Unless we say so explicitly, we target both non-degenerate and degenerate codes.

To begin, we display an example of a typical training trajectory in Fig. 2 to show the efficiency of our implementation. There, 16 agents are tasked to find $[[7, 1, 3]]$ codes, which each of them completes successfully running in parallel in 16 sec on a single GPU. The average circuit size starts being 25 by design, i.e. if no code has been found after 25 gates, the circuit gets reinitialized.

This number starts decreasing when codes start being found and it saturates to a final value, which is in general different for each agent. The noticeable variance in the circuit sizes that are found in Fig. 2 is a testament that different codes (in the sense of possibly belonging to different code families) with the same code parameters $[[7, 1, 3]]$ may have been found, yet making such a distinction requires further postprocessing.

We now move on to the main results of this Section, which are summarized in Fig. 3. There we show a classification of stabilizer codes with a range of different parameters $[[n, k, d]]$. In the following discussion we separate $d = 3, 4$ from $d = 5$, since the latter are more challenging to find.

For $d = 3$ and $d = 4$ codes we proceed as follows: for any given target $[[n, k, d]]$, we launch a few training runs. Once the codes are collected, we categorize them by calculating their quantum weight enumerators, see Eq. (8), leading to a certain number of non-degenerate ($A_{d-1} = B_{d-1} = 0$) and degenerate ($A_{d-1} = B_{d-1} \neq 0$) families. We repeat this process and keep launching new training runs until no new families are found by further runs. In this way, our strategy presumably finds *all* stabilizer codes that are possible for the given parameters n, k, d . This total number of families is shown in Fig. 3, with labels (x, y) for each $[[n, k, d]]$, where x is the number of non-degenerate families and y is the number of degenerate ones. It should be stressed that categorizing all stabilizer code families is in general an NP-complete problem [47], yet our framework is very effective at solving this task. To the best of our knowledge, this work provides the most detailed tabulation of (x, y) populations together with encoding circuits for the code parameters that we have studied.

This approach discovers suitable encoding circuits, given the assumed gate set, for a large set of codes. Among them are the following known codes for $d = 3$ (see [48] for explicit constructions of codes $[[n, n - r, 3]]$ with minimal r , for all n): The first one is the five-qubit perfect code [13], which consists of a *single* non-degenerate $[[5, 1, 3]]$ code family and is the smallest stabilizer code that corrects an arbitrary single-qubit error. Next are the 10 families [47] of $[[7, 1, 3]]$ codes, one of which corresponds to Steane’s code, the smallest $d = 3$ CSS code [14].

The smallest single-error-correcting surface code, Shor’s code [12], is rediscovered as one of the 143 degenerate code families with parameters $[[9, 1, 3]]$. The smallest quantum Hamming code [49] $[[8, 3, 3]]$ is obtained as well. Our approach is efficient enough to reach up to 20 physical qubits. The largest code parameters that we have considered for $d = 3$ are $[[20, 13, 3]]$, finding codes and encoding circuits in 800 sec with 45 gates (see Appendix B).

The circuit size shown for each $[[n, k, d]]$ in Fig. 3 is the minimal one found across all discovered families. In general, different families have different circuit sizes, and even within the same family we find variations in circuit sizes.

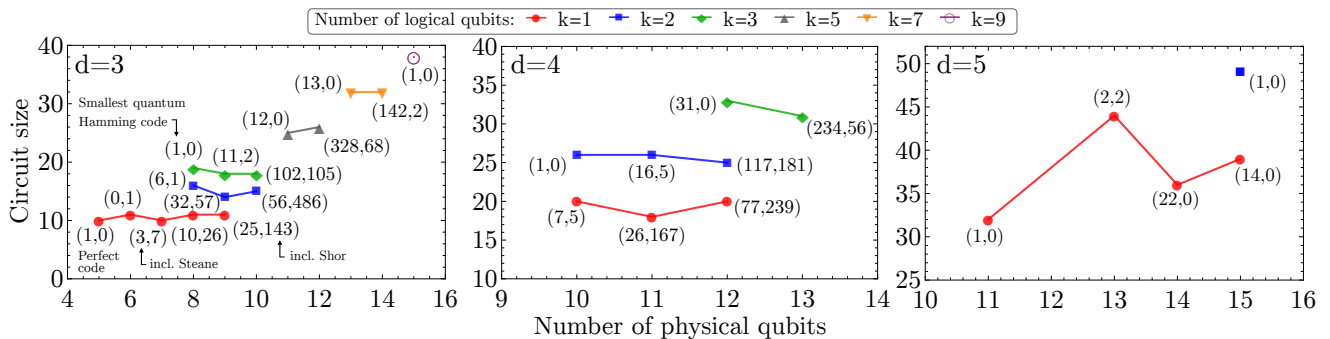


FIG. 3. Discovering codes and encoding circuits for various numbers of physical qubits, logical qubits, and distances. Selection of families of stabilizer codes tailored to symmetric depolarizing noise channels, found with our RL framework. The labels (x, y) indicate the number of non-degenerate (x) and degenerate (y) code families. The circuit size shown is the absolute minimum throughout all families, and different families in general have different minimal circuit sizes. Since further $d = 3, 4$ training runs do not increase family populations, it is likely that there are no more stabilizer codes for the shown $[[n, k, d = 3, 4]]$. Codes $[[n, k = 4, 6, 8, d = 3]]$ are also found but not shown to keep the figure uncluttered. For the computationally harder case of $d = 5$ we display the results found in a finite allotted time.

The RL framework presented here easily allows to find encoding circuits for different connectivities. The connectivity affects the likelihood of discovering codes within a certain family during RL training as well as the typical circuit sizes. In Fig. 4, we illustrate this for the case of $[[9, 3, 3]]$ codes, with their 13 families, for two different connectivities: an all-to-all (directed, i.e. $\text{CNOT}(i < j)$) and a nearest-neighbor square lattice connectivity.

We now move to distance $d = 5$ codes. These are more challenging to find due to the significantly increased number of error operators (11) to keep track of, which impacts both the computation time and the hardness of satisfying all KL conditions simultaneously. Nevertheless, our strategy is also successful in this case. It is known that the smallest possible distance-5 code has parameters $[[11, 1, 5]]$, a result that we confirm with our strategy. We find the single family for this code to have weight enumerators

$$A = (1, 0, 0, 0, 0, 0, 198, 0, 495, 0, 330, 0) , \quad (12)$$

$$B = (1, 0, 0, 0, 0, 0, 198, 198, 990, 495, 1650, 330, 234) ,$$

and an encoding circuit consisting of 32 gates in the minimal example, which we show in Appendix B.

In order to reduce computational effort, for $n \geq 14$ we ignored (5), and as a result the codes found in Fig. 3 $n \geq 14$ are only non-degenerate.

Moreover, the increased memory requirements from keeping track of more error operators (11) means that the number of agents that can be trained in parallel on a single GPU decreases. For instance, from the 4 agents that we train in parallel, it is rare that any of them finds an encoding sequence that leads to $d = 5$ code discovery. In addition, each of these training runs needs 1-4 hours, depending on the code parameters and whether degenerate codes are also targeted. Nevertheless, future performance improvements are likely possible.

B. Noise-aware meta-agent

We now move on to codes in more general asymmetric depolarizing noise channels. This lets us illustrate a powerful aspect of RL-based encoding and code discovery: One and the same agent can learn to switch its encoding strategy depending on some parameter characterizing the noise channel. This can be realized by training this noise-aware agent on many different runs with varying choices of the parameter, which is fed as an additional input to the agent. One motivation for this approach is that the agent may learn to generalize, i.e. transfer what it has learned between different values of the parameter.

In the present example, the parameter in question is the bias parameter $c_Z = \log p_Z / \log p_X$ introduced above, Eq. (6). This allows the *same* agent to switch its strategy depending on the kind of bias present in the noise channel. Once a particular value of c_Z is chosen, the error probabilities characterizing the noise channel are $(p_I, p_X, p_X, p_X^{c_Z})$. Normalization of the error probabilities imposes a relationship between p_I and p_X , which means that there is only one other free parameter besides c_Z , either p_I or p_X . It is more beneficial for training and generalization to keep p_I fixed and solve for p_X ; otherwise the magnitude of the probabilities $\{p_\mu\}$ changes a lot when varying c_Z , leading to poorer performance.

The error set E_μ is now taken to be all Pauli strings of (symmetric) weight ≤ 4 , i.e. $\{E_\mu\} = \{I, X_i, Y_i, Z_i, X_i X_j, \dots, Z_i Z_j Z_k Z_l\}$, but their associated error probabilities (and thus their effective weights according to (7)) will vary depending on c_Z . For every RL training trajectory, a new c_Z is chosen and the error probabilities p_μ are updated correspondingly. For the number of physical qubits that we will consider, the KL conditions (4), (5) cannot be exactly satisfied. Hence, we are forcing the agent to achieve some compromise: the most likely errors will have to be detected at the expense of not detecting other, less likely ones.

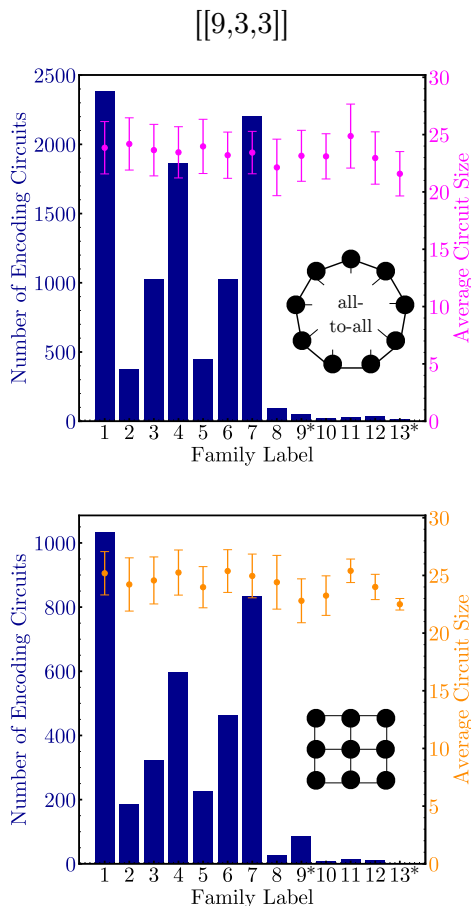


FIG. 4. Influence of connectivity. Characteristics of the 13 families of $[[9,3,3]]$ codes found with our framework, clustered according to families distinguished by their quantum weight enumerators (8). Families 9 and 13 (*) are degenerate, while the rest are non-degenerate. We have trained a total of 10240 agents for each of both cases. In the all-to-all (directed: $\text{CNOT}(i < j)$) connectivity, 9574 agents were successful, while this number went down to 3808 in the other case. The bars display how these codes are distributed across different families. Codes in the same family found by different agents are not necessarily distinct, so the bars are rather an indication of the likelihood of a training run to find a code within the family. The points show the mean circuit size, averaged within each family, while the error bar is its standard deviation. It is interesting to see that even with different connectivities, families occur with similar likelihoods during training. We explicitly list the corresponding quantum weight enumerators computed with (8) in Appendix A.

Regarding more detailed aspects of our implementation, we sample c_Z from the set $c_Z \in \{0.5, 0.6, 0.7, \dots, 1.9, 2\}$ with a uniform probability distribution. The hyperparameters λ_μ of the reward (10) are defined as

$$\lambda_\mu = \frac{p_\mu}{\max(p_\mu)} \Big|_{c_Z}, \quad (13)$$

by which we mean that for every c_Z , the corresponding set of p_μ 's gets normalized by the maximal value of p_μ in that set.

We choose $p_I = 0.9$, even though both slightly smaller and larger values around $p_I \approx 0.9$ perform equally well. However, going below $p_I \lesssim 0.8$ or above $p_I \gtrsim 0.95$ comes with different challenges. In the former (for large errors), we lose the important property that the sum of p_μ 's decreases as a function of weight, $(\sum_\mu p_\mu)_{w=1} > (\sum_\mu p_\mu)_{w=2} > \dots$. In the latter (small errors), the range of values of p_μ is so large that one would need to use a 64-bit floating-point representation to compute the reward with sufficient precision. Since both RL algorithms and GPUs are currently designed to work best with 32-bit precision, we decide to avoid this range of values for p_I during training, but we will still evaluate the strategies found by the RL agent on different values of p_I .

We apply this strategy to target codes with parameters $n = 9$, $k = 1$ in asymmetric noise channels. We allow a maximum number of 35 gates. Moreover, we consider an all-to-all connectivity, taking as available gate set $\{H_i, S_i, \text{CNOT}(i, j)\}$, where S_i is the phase gate acting on qubit i .

As is the case for most RL learning procedures, every independent learning run will typically result in a different learned strategy by the agent. We will thus train many agents and post-select the few best performing ones. Now, there are in principle two different ways to make this selection: The first one is based on how well they minimize the weighted KL sum (which is what they were trained for). The second one is by evaluating the probability that a single error correction cycle will end in failure, i.e. the probability that the wrong correction would be applied based on the detected syndrome. More concretely, we classify the agents by summing these two quantities over c_Z . By evaluating them based on these two criteria we will be able to see whether agents trained with a computationally cheaper reward (the weighted KL sum) can be reused for the more complex task of minimizing the failure probability. To be precise, we define the failure probability as

$$p_f = \sum_{\text{syn}} p(\text{correct wrong error}|\text{syn})p(\text{syn}), \quad (14)$$

where syn is one of the possible 2^{n-k} syndromes, and where the error correction strategy is to always correct the most likely error to have happened given that specific syndrome. If more than one error is equally likely, we randomly select one of them.

We discover codes with the following parameters: $[[9, 1, d_e(c_Z = 0.5) = 2]]$, $[[9, 1, d_e(c_Z = 0.6) = 3]]$, $[[9, 1, d_e(c_Z = 1.4) = 4]]$, $[[9, 1, d_e(c_Z = 2) = 5]]$. Codes in-between, $0.5 \leq c_Z < 0.6$, have $d_e = 2$, $0.6 \leq c_Z < 1.4$ have $d_e = 3$, and so on. Even though all encodings that the agent outputs have circuit size 35, we notice that trivial gate sequences are applied at the last few steps, effectively reducing the overall gate count. We remark that this feature is not problematic: it means that the

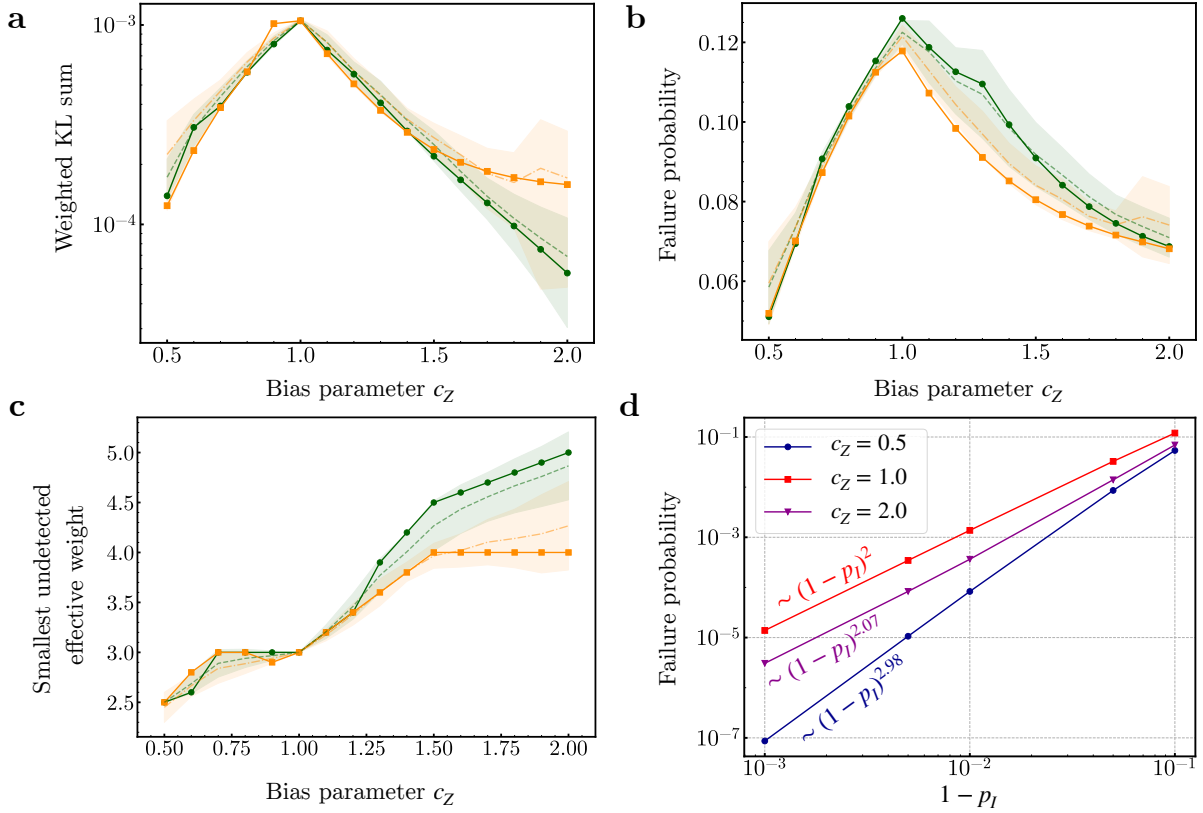


FIG. 5. Performance of the noise-aware RL agent. The agent finds $n = 9, k = 1$ codes and encoding circuits, simultaneously for different levels of noise bias c_Z , with single-qubit fidelity $p_I = 0.9$. In panels **a,b,c**, green represents the agent that was post-selected among all trained agents for performing best at minimizing the weighted KL sum, averaged over all c_Z values. Orange refers to the agent minimizing the failure probability, averaged over c_Z . Dashed lines correspond to averages and shaded regions to the standard deviations over the best 30 agents of each class. **a** Weighted KL sum as a function of the noise bias parameter c_Z (best agent: solid green line). **b** Failure probability as a function of the noise bias parameter c_Z . Here the best agent, as measured in this way, follows the solid orange line with square markers. **c** Smallest undetected effective weight (effective code distance is the integer part) as a function of the noise bias parameter c_Z . While there is almost a perfect overlap between both ‘best’ agents until $c_Z = 1.1$, the situation changes afterwards, leading at $c_Z = 2$ to a $d_e = 5$ code (green) or a $d_e = 4$ code (orange) that perform equally well in terms of the failure probability, as seen in **b**. **d** Evaluation of the failure-probability of the best-performing agent (orange in the other panels) for larger values of p_I (smaller errors) than the ones it was trained on.

agent is done well before a new training run is launched, and the best thing it can do is collecting small negative rewards until the end. We manually prune the encodings to get rid of such trivial operations, and the resulting circuit sizes vary from 22 to 35, depending on the value of c_Z .

The main results are shown in Fig. 5. We start by comparing the two best-performing post-selected agents according to minimizing the weighted KL sum (green) and minimizing the failure probability (orange), see Fig. 5a,b. There we see that it is in general not true that minimizing the weighted KL sum will automatically correspond to a lower failure probability. Nevertheless, there is a nice correlation between the two, especially in the region $c_Z < 1$. More explicitly, the two best-performing agents show an almost identical performance in terms of both the weighted KL sum and the failure probability in this region. The situation changes when $c_Z \geq 1$, where we see

that the optimal agents for the weighted KL sum are not optimal in terms of minimizing the failure probability, especially when $c_Z \geq 1.5$.

The results from Fig. 5a,b suggest that there are two “branches” separated by a “phase transition” at $c_Z = 1$. Indeed, the two-branch structure can be easily understood in terms of the degree of asymmetry in the noise channel. For $c_Z < 1$, Z -errors are more likely than X/Y -errors, while the opposite happens for $c_Z > 1$. For $c_Z < 1$, a single error “species” had to be prioritized, but this number grows to two for $c_Z > 1$. At the middle point, $c_Z = 1$, there is no preference for any of the three species.

The last comparison between these two agents that we do is based on the smallest undetected effective weight of the codes found in Fig. 5c. The same behavior as before is observed: almost-perfect correlation up to $c_Z = 1$, and discrepancy at larger values of c_Z . Surprisingly, the code

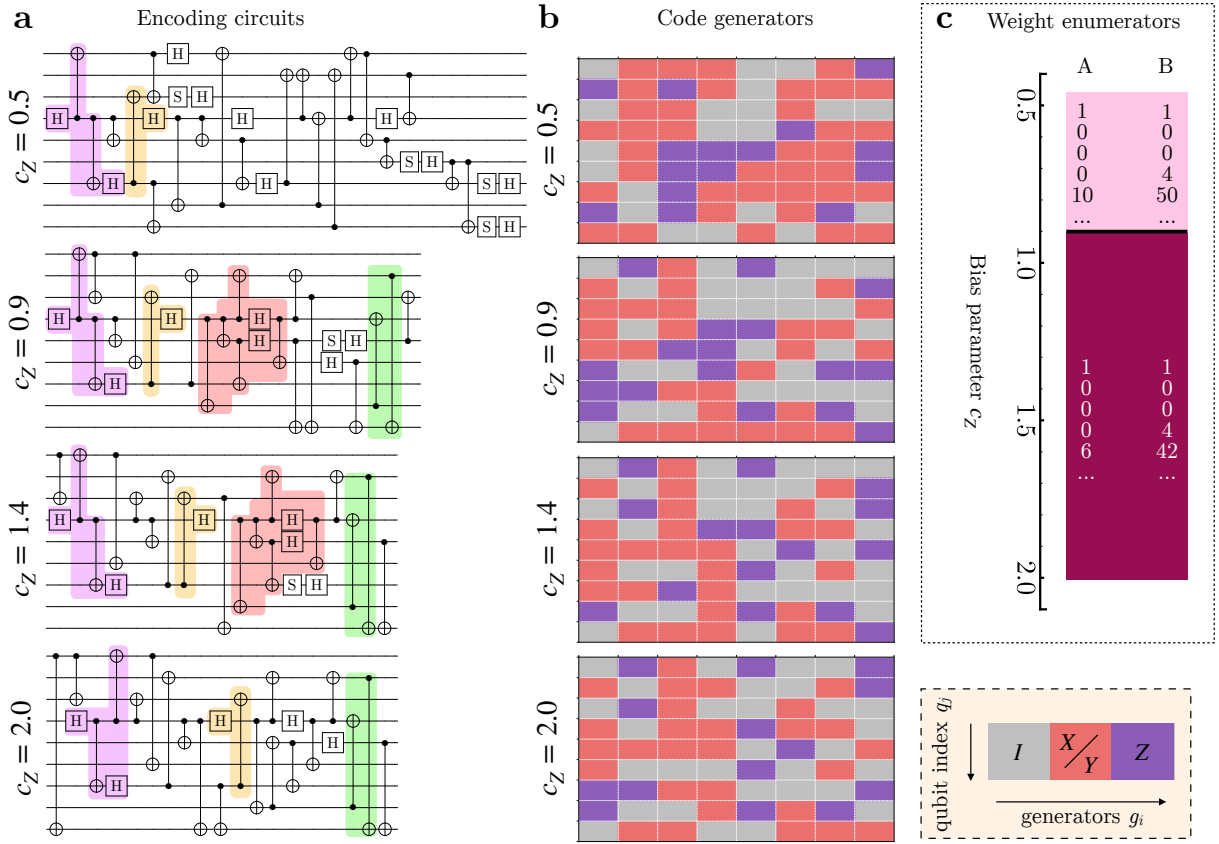


FIG. 6. Characteristics of the 9-qubit codes and encodings found by the noise-aware meta-agent post-selected for minimizing the failure probability. **a** Encoding circuits: Here we see that many small gate sequences (highlighted with different colors) are reused across different values of c_Z . This is an indication of transfer learning, i.e. the power of the meta-agent. **b** Code generators g_i corresponding to the encoding circuits, where we do not make a distinction between X or Y . Here we see that the code generators g_i vary across different values of c_Z . **c** Associated code family according to their (symmetric) weight enumerators A, B . The same code family is used from $0.5 \leq c_Z < 0.9$, while a family switching occurs at $c_Z = 0.9$, and it is kept until $c_Z = 2$.

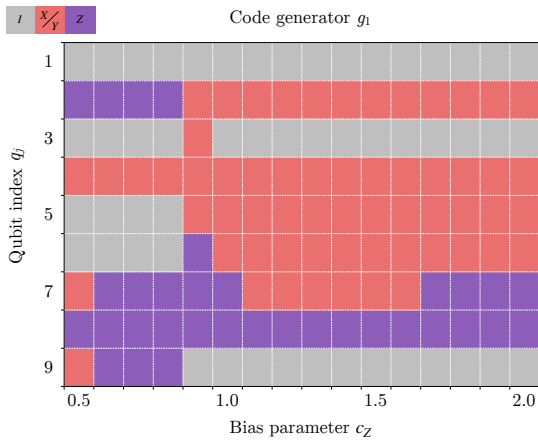


FIG. 7. Evolution of the first code generator g_1 as a function of the bias parameter c_Z .

found by the best agent according to the weighted KL

sum (green) at $c_Z = 2$ has $d_e = 5$, while the best code at minimizing the failure probability (orange) has $d_e = 4$. However, at the specific point $c_Z = 2$ these two codes perform equally in terms of the failure probability (see Fig. 5b).

Now, we focus on the agent that performs best at minimizing the failure probability (orange) since it is the one of most interest in physical scenarios. We are specifically interested in doing a more detailed analysis to understand what are the characteristics of the codes that are being found. We begin by evaluating the performance of the *same* agent on different values of p_I . This is shown in Fig. 5d, where we see that there is indeed a gain by encoding the quantum state, and it asymptotically follows a power law with exponent $\gtrsim 2$ depending on the specific value of c_Z .

We continue by analyzing the encoding circuits and code generators for some selected values of c_Z . These are chosen after computing the (symmetric) quantum weight enumerators according to Eq. (8), which we show

in Fig. 6c. There we see that the same code family is kept for $0.5 \leq c_Z < 0.9$, where Z errors are more likely than X/Y. From that point onward, the agent switches to a new code family that is kept until the end, $c_Z = 2$. We thus choose to analyze the encoding circuits and their associated code generators for the values $c_Z = \{0.5, 0.9, 1.4, 2\}$.

We begin by showing the encoding circuits in Fig. 6a, highlighting common motifs that are re-used across various values of c_Z with different colors, indicative of transfer learning. Another interesting behavior is that S gates are used more prominently at small values of c_Z , in particular in the combination $S \cdot H$. What this combination of gates does is a permutation: $X \rightarrow Y, Y \rightarrow Z, Z \rightarrow X$ (ignoring signs), which is very useful to exchange Y by Z efficiently. Moreover, we recall that the circuits that we show are pruned, i.e. trivial sequences of gates are manually eliminated. The agent actually uses lots of S gates at the very end of the encoding for almost all values of c_Z . This is no accident: we have constructed a scenario in which X/Y errors are essentially indistinguishable, in the sense that they occur with equal probability. Thus the S gate (which exchanges X by Y) acts like an idling operation. This is cleverly leveraged by the agent to shorten the encoding circuits: once it cannot improve the final reward further, the strategy then becomes to reach that point as soon as possible and collect that final reward for as many steps as possible.

The final aspect we investigate is the code generators of such encoding circuits, which are shown in Fig. 6b. To aid visualization, we have chosen different colors for different Pauli matrices. However, since our scenario is by construction symmetric in X/Y, we choose to represent X and Y by the same color. Since the code used at $c_Z = 0.5$ is the only one from a different code family, it is natural that its code generator pattern is the most distinct. However, we see that the generators of the remaining values of c_Z have similar structures at certain points. We also include in Fig. 7 how the first code generator g_1 (the first column in Fig. 6) changes across all values of c_Z . For instance, the transition between different code families is clearly visible at $c_Z = 0.9$. In addition, other

(not shown) features such as the agent using the *same* circuit throughout $1.2 \leq c_Z \leq 1.6$ can also be glimpsed from Fig. 7 (even though one would need to analyze how all the generators evolve).

V. CONCLUSIONS AND OUTLOOK

We have presented an efficient RL framework that is able to simultaneously discover QEC codes and their encoding circuits from scratch, given a qubit connectivity, gate set, and error operators. It learns strategies simultaneously for a range of noise models, thus re-using and transferring discoveries between different noise regimes. We have been able to discover codes and circuits up to 20 physical qubits and code distance 5. This is thanks to our formulation in terms of stabilizers, that serve both as compact input to the agent as well as the basis for rapid Clifford simulations, which we implemented in a vectorized fashion using a modern machine-learning framework.

In the present work, we have focused on the quantum communication or quantum memory scenario, where the encoding circuit itself can be assumed error-free since we focus on errors happening during transmission. As a result, our encoding circuits are not fault tolerant, i.e. single errors, when introduced, might sometimes proliferate to become incorrigible. Flag-based fault tolerance [50] added on top of our encoding circuits could turn them fault tolerant. Another very interesting extension of our scheme would be to design rewards that directly target codes that admit transverse logical gates or have other useful aspects, such as the requirements in [16] for codes that can outperform the surface code.

ACKNOWLEDGMENTS

Fruitful discussions with Sangkha Borah, Jonas Landgraf and Maximilian Naegele are thankfully acknowledged. This research is part of the Munich Quantum Valley, which is supported by the Bavarian state government with funds from the Hightech Agenda Bayern Plus.

-
- [1] A. Acín, I. Bloch, H. Buhrman, T. Calarco, C. Eichler, J. Eisert, D. Esteve, N. Gisin, S. J. Glaser, F. Jelezko, S. Kuhr, M. Lewenstein, M. F. Riedel, P. O. Schmidt, R. Thew, A. Wallraff, I. Walmsley, and F. K. Wilhelm, The quantum technologies roadmap: a european community view, *New Journal of Physics* **20**, 080201 (2018).
 - [2] S. M. Girvin, Introduction to quantum error correction and fault tolerance, *SciPost Physics Lecture Notes* 10.21468/scipostphyslectnotes.70 (2023).
 - [3] M. Inguscio, W. Ketterle, and C. Salomon, *Proceedings of the International School of Physics "Enrico Fermi."*, Vol. 164 (IOS press, 2007).
 - [4] J. Preskill, Fault-tolerant quantum computation (1997), arXiv:quant-ph/9712048 [quant-ph].
 - [5] A. Paler and S. J. Devitt, An introduction into fault-tolerant quantum computing, in *2015 52nd ACM/EDAC/IEEE Design Automation Conference (DAC)* (2015) pp. 1–6.
 - [6] S. Krinner, N. Lacroix, A. Remm, A. Di Paolo, E. Genois, C. Leroux, C. Hellings, S. Lazar, F. Swiadek, J. Herrmann, *et al.*, Realizing repeated quantum error correction in a distance-three surface code, *Nature* **605**, 669 (2022).
 - [7] C. Ryan-Anderson, J. G. Bohnet, K. Lee, D. Gresh, A. Hankin, J. Gaebler, D. Francois, A. Chernoguzov, D. Lucchetti, N. C. Brown, *et al.*, Realization of real-time

- fault-tolerant quantum error correction, *Physical Review X* **11**, 041058 (2021).
- [8] L. Postler, S. Heußen, I. Pogorelov, M. Rispler, T. Feldker, M. Meth, C. D. Marciniak, R. Stricker, M. Ringbauer, R. Blatt, *et al.*, Demonstration of fault-tolerant universal quantum gate operations, *Nature* **605**, 675 (2022).
- [9] I. Cong, H. Levine, A. Keesling, D. Bluvstein, S.-T. Wang, and M. D. Lukin, Hardware-efficient, fault-tolerant quantum computation with rydberg atoms, *Physical Review X* **12**, 021049 (2022).
- [10] R. Acharya, I. Aleiner, R. Allen, *et al.*, Suppressing quantum errors by scaling a surface code logical qubit, *Nature* **614**, 676 (2023).
- [11] V. Sivak, A. Eickbusch, B. Royer, S. Singh, I. Tsioutsios, S. Ganjam, A. Miano, B. Brock, A. Ding, L. Frunzio, *et al.*, Real-time quantum error correction beyond break-even, *Nature* **616**, 50 (2023).
- [12] A. R. Calderbank and P. W. Shor, Good quantum error-correcting codes exist, *Phys. Rev. A* **54**, 1098 (1996).
- [13] R. Laflamme, C. Miquel, J. P. Paz, and W. H. Zurek, Perfect quantum error correcting code, *Phys. Rev. Lett.* **77**, 198 (1996).
- [14] A. M. Steane, Simple quantum error-correcting codes, *Phys. Rev. A* **54**, 4741 (1996).
- [15] A. Y. Kitaev, Quantum computations: algorithms and error correction, *Russian Mathematical Surveys* **52**, 1191 (1997).
- [16] S. Bravyi, A. W. Cross, J. M. Gambetta, D. Maslov, P. Rall, and T. J. Yoder, High-threshold and low-overhead fault-tolerant quantum memory (2023), arXiv:2308.07915 [quant-ph].
- [17] M. Grassl and S. Han, Computing extensions of linear codes using a greedy algorithm, in *2012 IEEE International Symposium on Information Theory Proceedings* (2012) pp. 1568–1572.
- [18] M. Grassl, P. W. Shor, G. Smith, J. Smolin, and B. Zeng, New constructions of codes for asymmetric channels via concatenation, *IEEE Transactions on Information Theory* **61**, 1879 (2015).
- [19] M. Li, M. Gutiérrez, S. E. David, A. Hernandez, and K. R. Brown, Fault tolerance with bare ancillary qubits for a $[[7,1,3]]$ code, *Phys. Rev. A* 10.1103/PhysRevA.96.032341 (2017).
- [20] I. Chuang, A. Cross, G. Smith, J. Smolin, and B. Zeng, Codeword stabilized quantum codes: Algorithm and structure, *Journal of Mathematical Physics* 10.1063/1.3086833 (2009), https://pubs.aip.org/aip/jmp/article-pdf/doi/10.1063/1.3086833/15610997/042109_1_online.pdf.
- [21] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, Policy gradient methods for reinforcement learning with function approximation, *Advances in neural information processing systems* **12** (1999).
- [22] T. Fösel, P. Tighineanu, T. Weiss, and F. Marquardt, Reinforcement learning with neural networks for quantum feedback, *Phys. Rev. X* 10.1103/PhysRevX.8.031084 (2018).
- [23] H. P. Nautrup, N. Delfosse, V. Dunjko, H. J. Briegel, and N. Friis, Optimizing quantum error correction codes with reinforcement learning, *Quantum* 10.22331/q-2019-12-16-215 (2019).
- [24] C. Maun, T. Farrelly, and T. M. Stace, Optimization of tensor network codes with reinforcement learning, arXiv:2305.11470 (2023).
- [25] V. P. Su, C. Cao, H.-Y. Hu, Y. Yanay, C. Tahan, and B. Swingle, Discovery of optimal quantum error correcting codes via reinforcement learning (2023), arXiv:2305.06378 [quant-ph].
- [26] C. Cao and B. Lackey, Quantum lego: Building quantum error correction codes from tensor networks, *PRX Quantum* **3**, 020332 (2022).
- [27] P. Andreasson, J. Johansson, S. Liljestrand, and M. Granath, Quantum error correction for the toric code using deep reinforcement learning, *Quantum* **3**, 183 (2019).
- [28] R. Sweke, M. S. Kesselring, E. P. van Nieuwenburg, and J. Eisert, Reinforcement learning decoders for fault-tolerant quantum computation, *Machine Learning: Science and Technology* **2**, 025005 (2020).
- [29] L. D. Colomer, M. Skotiniotis, and R. Muñoz-Tapia, Reinforcement learning for optimal error correction of toric codes, *Physics Letters A* **384**, 126353 (2020).
- [30] C. Cao, C. Zhang, Z. Wu, M. Grassl, and B. Zeng, Quantum variational learning for quantum error-correcting codes, *Quantum* 10.22331/q-2022-10-06-828 (2022).
- [31] D. Gottesman, Stabilizer codes and quantum error correction (1997), quant-ph/9705052 [quant-ph].
- [32] M. A. Nielsen and I. L. Chuang, *Quantum computation and quantum information* (Cambridge university press, 2010).
- [33] C. H. Bennett, D. P. DiVincenzo, J. A. Smolin, and W. K. Wootters, Mixed-state entanglement and quantum error correction, *Phys. Rev. A* **54**, 3824 (1996).
- [34] E. Knill and R. Laflamme, Theory of quantum error-correcting codes, *Phys. Rev. A* 10.1103/PhysRevA.55.900 (1997).
- [35] L. Ioffe and M. Mézard, Asymmetric quantum error-correcting codes, *Phys. Rev. A* 10.1103/PhysRevA.75.032345 (2007).
- [36] L. Wang, K. Feng, S. Ling, and C. Xing, Asymmetric quantum codes: Characterization and constructions, *IEEE Transactions on Information Theory* **56**, 2938 (2010).
- [37] M. F. Ezerman, S. Ling, and P. Sole, Additive asymmetric quantum codes, *IEEE Transactions on Information Theory* **57**, 5536 (2011).
- [38] G. G. L. Guardia, On the construction of asymmetric quantum codes, *International Journal of Theoretical Physics* 10.1007/s10773-014-2031-y (2014).
- [39] P. Shor and R. Laflamme, Quantum analog of the macwilliams identities for classical coding theory (1997).
- [40] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction* (MIT press, 2018).
- [41] V. Konda and J. Tsitsiklis, Actor-critic algorithms, *Advances in neural information processing systems* **12** (1999).
- [42] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal policy optimization algorithms, arXiv:1707.06347 (2017).
- [43] C. Gidney, Stim: a fast stabilizer circuit simulator, *Quantum* **5**, 497 (2021).
- [44] S. Aaronson and D. Gottesman, Improved simulation of stabilizer circuits, *Phys. Rev. A* 10.1103/physreva.70.052328 (2004).
- [45] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, A. P. George Necula, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, JAX: com-

posable transformations of Python+NumPy programs (2018).

- [46] C. Lu, J. Kuba, A. Letcher, L. Metz, C. S. de Witt, and J. Foerster, Discovered policy optimisation, *Advances in Neural Information Processing Systems* , 16455 (2022).
- [47] S. Yu, Q. Chen, and C. H. Oh, Graphical quantum error-correcting codes (2007), arXiv:0709.1780 [quant-ph].
- [48] S. Yu, J. Bierbrauer, Y. Dong, Q. Chen, and C. Oh, All the stabilizer codes of distance 3, *IEEE transactions on information theory* **59**, 5179 (2013).
- [49] D. Gottesman, Class of quantum error-correcting codes saturating the quantum hamming bound, *Phys. Rev. A* **54**, 1862 (1996).
- [50] R. Chao and B. W. Reichardt, Quantum error correction with only two extra qubits, *Phys. Rev. Lett.* 10.1103/physrevlett.121.050502 (2018).

Appendix A: Quantum Weight Enumerators of $[[9, 3, 3]]$

Here we list the quantum weight enumerators of the 13 families of $[[9, 3, 3]]$ codes appearing in the main text, ordered accordingly.

$$\begin{aligned} A &= (1, 0, 0, 0, 0, 6, 16, 24, 15, 2) , \\ B &= (1, 0, 0, 40, 162, 480, 952, 1224, 933, 304) , \end{aligned} \quad (\text{A1})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 2, 4, 12, 28, 17, 0) , \\ B &= (1, 0, 0, 48, 146, 472, 984, 1216, 917, 312) , \end{aligned} \quad (\text{A2})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 2, 2, 16, 28, 13, 2) , \\ B &= (1, 0, 0, 44, 162, 452, 984, 1236, 901, 316) , \end{aligned} \quad (\text{A3})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 1, 3, 18, 26, 12, 3) , \\ B &= (1, 0, 0, 40, 170, 456, 968, 1240, 909, 312) , \end{aligned} \quad (\text{A4})$$

$$\begin{aligned} A &= (1, 0, 0, 1, 0, 3, 16, 27, 15, 1) , \\ B &= (1, 0, 0, 48, 162, 456, 952, 1248, 933, 296) , \end{aligned} \quad (\text{A5})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 0, 4, 20, 24, 11, 4) , \\ B &= (1, 0, 0, 36, 178, 460, 952, 1244, 917, 308) , \end{aligned} \quad (\text{A6})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 1, 5, 14, 26, 16, 1) , \\ B &= (1, 0, 0, 44, 154, 476, 968, 1220, 925, 308) , \end{aligned} \quad (\text{A7})$$

$$\begin{aligned} A &= (1, 0, 0, 1, 0, 1, 20, 27, 11, 3) , \\ B &= (1, 0, 0, 44, 178, 436, 952, 1268, 917, 300) , \end{aligned} \quad (\text{A8})$$

$$\begin{aligned} A &= (1, 0, 1, 0, 0, 4, 13, 28, 17, 0) , \\ B &= (1, 0, 1, 57, 171, 481, 931, 1171, 944, 339) , \end{aligned} \quad (\text{A9})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 4, 0, 12, 32, 15, 0) , \\ B &= (1, 0, 0, 52, 146, 444, 1016, 1228, 885, 324) , \end{aligned} \quad (\text{A10})$$

$$\begin{aligned} A &= (1, 0, 0, 0, 0, 0, 36, 0, 27, 0) , \\ B &= (1, 0, 0, 36, 162, 540, 792, 1404, 837, 324) , \end{aligned} \quad (\text{A11})$$

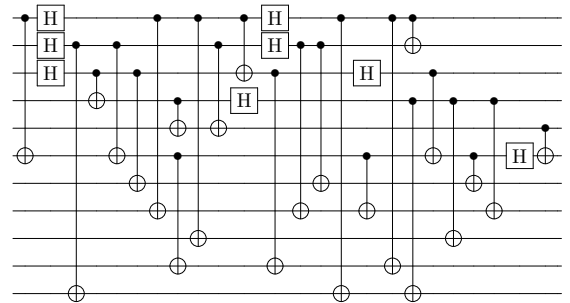
$$\begin{aligned} A &= (1, 0, 0, 2, 0, 0, 16, 30, 15, 0) , \\ B &= (1, 0, 0, 56, 162, 432, 952, 1272, 933, 288) , \end{aligned} \quad (\text{A12})$$

$$\begin{aligned} A &= (1, 0, 1, 0, 0, 0, 21, 28, 9, 4) , \\ B &= (1, 0, 1, 49, 203, 441, 931, 1211, 912, 347) , \end{aligned} \quad (\text{A13})$$

Appendix B: Encoding circuits

Here we show some minimal encoding circuits for selected code parameters.

$[[11, 1, 5]]$: 32 gates in a directed all-to-all connectivity:



$[[20, 13, 3]]$: 45 gates in an all-to-all connectivity:

