

# Reinforcement Learning with Neural Networks for Quantum Feedback

Thomas Fösel, Petru Tighineanu, and Talitha Weiss

*Max Planck Institute for the Science of Light, Staudtstr. 2, 91058 Erlangen, Germany*

Florian Marquardt

*Max Planck Institute for the Science of Light, Staudtstr. 2, 91058 Erlangen, Germany and Physics Department, University of Erlangen-Nuremberg, Staudtstr. 5, 91058 Erlangen, Germany*

Machine learning with artificial neural networks is revolutionizing science. The most advanced challenges require discovering answers autonomously. This is the domain of reinforcement learning, where control strategies are improved according to a reward function. The power of neural-network-based reinforcement learning has been highlighted by spectacular recent successes, such as playing Go, but its benefits for physics are yet to be demonstrated. Here, we show how a network-based “agent” can discover complete quantum-error-correction strategies, protecting a collection of qubits against noise. These strategies require feedback adapted to measurement outcomes. Finding them from scratch, without human guidance, tailored to different hardware resources, is a formidable challenge due to the combinatorially large search space. To solve this, we develop two ideas: two-stage learning with teacher/student networks and a reward quantifying the capability to recover the quantum information stored in a multi-qubit system. Beyond its immediate impact on quantum computation, our work more generally demonstrates the promise of neural-network-based reinforcement learning in physics.

We are witnessing rapid progress in applications of artificial neural networks (ANN) for tasks like image classification, speech recognition, natural language processing, and many others [1, 2]. Within physics, the examples emerging during the past two years range across areas like statistical physics, quantum many-body systems, and quantum error correction [3–11]. To date, most applications of neural networks employ supervised learning, where a large collection of samples has to be provided together with the correct labeling.

However, inspired by the long-term vision of artificial scientific discovery [12, 13], one is led to search for more powerful techniques that explore solutions to a given task autonomously. Reinforcement learning (RL) is a general approach of this kind [2], where an “agent” interacts with an “environment”. The agent’s “policy”, i.e. the choice of actions in response to the environment’s evolution, is updated to increase some reward. The power of this method, when combined with ANNs, was demonstrated convincingly through learning to play games beyond human expertise [14, 15]. In physics, RL *without* neural networks has been introduced recently, for example to study qubit control [16, 17] and invent quantum optics experiments [18]. Moving to neural-network-based RL promises access to the vast variety of techniques currently being developed for ANNs.

In this work, we introduce network-based RL in physics (Fig. 1) and illustrate its versatility in the domain of quantum feedback. Specifically, we devise a unified, fully autonomous, human-guidance-free approach for discovering quantum-error-correction (QEC) strategies from scratch, in few-qubit quantum systems subject to arbitrary noise and hardware constraints. This approach relies on a network agent that learns feedback strategies, adapting its actions to measurement results. As illustrated in Fig. 1b-d, our method provides a unified approach to protect a

quantum memory from noise. It covers a wide range of scenarios where one would otherwise have to select an existing scheme (stabilizer codes, adaptive phase estimation, etc.) and adapt it to the given situation. Our findings are of immediate relevance to the broad field of quantum error correction (including quantum-error-mitigation techniques) and are best suited to be used in few-qubit quantum modules. These could be used as stand-alone quantum memory or be part of the modular approach to quantum computation, which has been suggested for several leading hardware platforms[19, 20].

Given a collection of qubits and a set of available quantum gates, the agent is asked to preserve an arbitrary quantum state  $\alpha|0\rangle + \beta|1\rangle$  initially stored in one of the qubits. It finds complex sequences including projective measurements and entangling gates, thereby protecting the quantum information stored in such a few-qubit system against decoherence. This is a very complex challenge, where both brute force searches and even the most straightforward RL approaches fail. The success of our approach is due to a combination of two key ideas: (i) ~~two-stage learning, with an RL-trained network receiving maximum input acting as a teacher for a second network, and (ii) a measure of the recoverable quantum information hidden inside a collection of qubits, being used as a reward.~~

Recent progress in multi-qubit quantum devices [21–30] has highlighted hardware features deviating from often-assumed idealized scenarios. These include qubit connectivity, correlated noise, restrictions on measurements, or inhomogeneous error rates. Our approach can help finding “hardware-adapted” solutions. This builds on the main advantage of RL, namely its flexibility: it can discover strategies for such a wide range of situations with minimal domain-specific input. We illustrate this flexibility in examples from two different domains: in

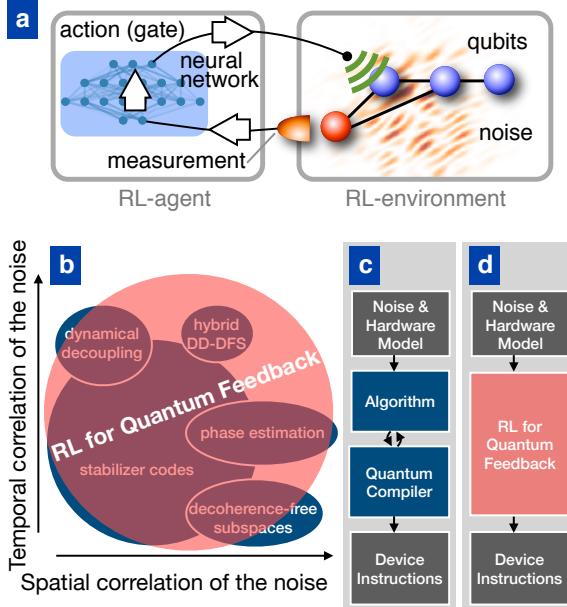


Figure 1. (color) (a) The general setting of this work: A few-qubit quantum device with a neural-network-based controller whose task is to protect the quantum memory residing in this device against noise. Reinforcement learning (RL) lets the controller (“RL-agent”) discover on its own how to best choose gate sequences, perform measurements, and react to measurement results, by interacting with the quantum device (“RL-environment”). (b) visualizes the flexibility of our approach (schematic). Depending on the type of noise and hardware setting, different approaches are optimal (DD, dynamical decoupling; DFS, decoherence-free subspace). By contrast, the RL approach is designed to automatically discover the best strategy, adapted to the situation. In (c) we show the conventional procedure to select some QEC algorithm and then produce hardware-adapted device instructions (possibly re-iterating until an optimal choice is found). We compare this to our approach (d) that takes care of all these steps at once and provides QEC strategies fully adapted to the concrete specifications of the quantum device.

one set of examples (uncorrelated bit-flip noise), the network is able to go beyond rediscovering the textbook stabilizer repetition code. It finds an adaptive response to unexpected measurement results that allows it to increase the coherence time, performing better than any straightforward non-adaptive implementation. Simultaneously, it automatically discovers suitable gate sequences for various types of hardware settings. In another, very different example, the agent learns to counter spatially correlated noise by finding non-trivial adaptive **phase-estimation strategies** that quickly become intractable by conventional numerical approaches such as brute-force search. Crucially, all these examples can be treated by exactly the same approach, with no fine-tuning. The only input consists in the problem specification (hardware and noise model).

In a nutshell, our goal is to have a neural network which

can be employed in an experiment, receiving measurement results and selecting suitable subsequent gate operations conditioned on these results. However, in our two-stage learning approach, we do not directly train this neural network from scratch. Rather, we first employ reinforcement learning to train an auxiliary network that has full knowledge of the simulated quantum evolution. Later on, the experimentally applicable network is trained in a supervised way to mimic the behavior of this auxiliary network.

We emphasize that feedback requires reaction towards the observations, going beyond optimal control type challenges (like pulse shape optimization or dynamical decoupling), and RL has been designed for exactly this purpose. Specifically, in this work we will consider discrete-time, digital feedback, of the type that is now starting to be implemented experimentally [31–35], e.g. for error correction in superconducting quantum computers. Other wide-spread optimization techniques for quantum control, like GRAPE, often vary evolution operators with respect to continuous parameters [36, 37], but do not easily include feedback and are most suited for optimizing the pulse shapes of individual gates (rather than complex gate sequences acting on many qubits). Another recent approach [38] to quantum error correction uses optimization of control parameters in a pre-configured gate sequence. By contrast, RL directly explores the space of discrete gate sequences. Moreover, it is a “model-free” approach [2], i.e. it does not rely on access to the underlying dynamics. What is optimized is the network agent. Neural-network based RL promises to complement other successful machine-learning techniques applied to quantum control [39–42].

Conceptually, our approach aims to control a quantum system using a classical neural network. To avoid confusion, we emphasize our approach is distinct from future “quantum machine learning” devices, where even the network will be quantum [8, 43, 44].

## Reinforcement Learning

The purpose of RL (Fig. 1a) is to find an optimal set of actions (in our case, quantum gates and measurements) that an “agent” can perform in response to the changing state of an “environment” (here, the quantum memory). The objective is to maximize the expected “return”  $R$ , i.e. a sum of rewards.

To find optimal gate sequences, we employ a widespread version of reinforcement learning [45, 46] where discrete actions are selected at each time step  $t$  according to a probabilistic “policy”  $\pi_\theta$ . Here,  $\pi_\theta(a_t|s_t)$  is the probability to apply action  $a_t$ , given the state  $s_t$  of the RL-environment. As we will use a neural network to compute  $\pi_\theta$ , the multi-dimensional parameter  $\theta$  stands for all the network’s weights and biases. The network is fed  $s_t$  as an input vector and outputs the probabilities  $\pi_\theta(a_t|s_t)$ . The expected return can then be maximized by applying the

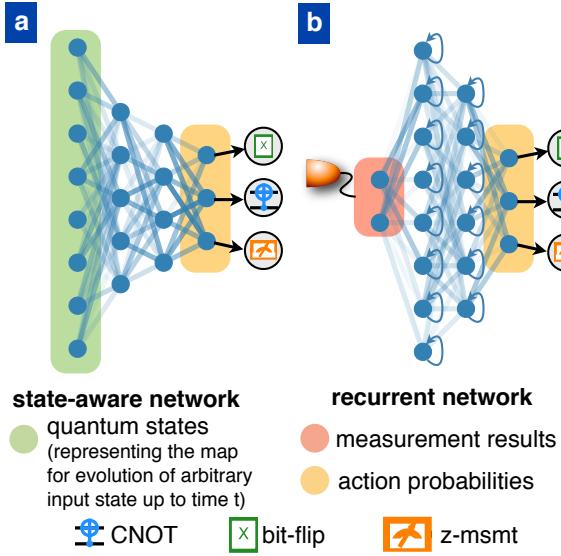


Figure 2. (color) The neural networks. (a) At each time  $t$ , the “state-aware” network receives a representation of the map  $\Phi$  describing the quantum evolution of arbitrary initial logical qubit states up to that time (represented by four evolved states  $\hat{\rho}$ ; see main text). It outputs the probabilities for the different actions (gates), defining the agent’s policy. Like any neural network, it is a nonlinear function that can be decomposed into layer-wise linear superposition of neuron values, using trainable weights (visualized by connections), and the application of nonlinear activation functions. Examples for actions are shown here (bit-flip, CNOT, measurement). Each of those can be applied to various qubits (or qubit pairs), resulting in around 10-20 actions. (b) The recurrent network receives the most recent measurement result (if any) and also outputs action probabilities. Its long short-term memory (LSTM) neurons learn to keep track of information accumulated in previous time steps (schematically indicated by the recurrent connections here).

policy gradient RL update rule [45]:

$$\delta\theta_j = \eta \frac{\partial \mathbb{E}[R]}{\partial \theta_j} = \eta \mathbb{E} \left[ R \sum_t \frac{\partial}{\partial \theta_j} \ln \pi_\theta(a_t | s_t) \right], \quad (1)$$

with  $\eta$  the learning rate parameter, and  $\mathbb{E}$  the expectation value over all gate sequences and measurement outcomes. These ingredients summarize the basic policy gradient approach. In practice, improvements of Eq. (1) are used; for example, we employ a baseline, natural policy gradient, and entropy regularization (see Appendix). Even so, several further conceptual steps are essential to have any chance of success (see below).

Eq. (1) provides the standard recipe for a fully observed environment. This approach can be extended to a partially observed environment, where the policy would then be a function of the observations only, instead of the state. The observations contain partial information on the actual state of the environment. In the present manuscript we will encounter both cases.

## Reinforcement Learning Approach to Quantum Memory

In this work we seek to train a neural network to develop strategies to protect the quantum information stored in a quantum memory from decoherence. This involves both variants of stabilizer-code-based QEC [47–49] as well as other, more specialized (but, in their respective domain, more resource-efficient) approaches, like decoherence-free subspaces or phase estimation. We remind the reader that, for the particular case of stabilizer-code-based QEC, the typical steps are: (i) the encoding, in which the logical state initially stored in one qubit is distributed over several physical qubits, (ii) the detection of errors via measurement of suitable multi-qubit operators (syndromes), (iii) the subsequent correction, and (iv) the decoding procedure that transfers the encoded state back into one physical qubit. We stress that no such specialized knowledge will be provided a priori to our network, thus retaining maximum flexibility in the tasks it might be applied to and in the strategies it can encompass (Fig. 1b).

We start by storing an arbitrary quantum state  $\alpha|0\rangle + \beta|1\rangle$  inside one physical qubit. The goal is to be able to retrieve this state with optimum fidelity after a given time span. Given hardware constraints such as the connectivity between qubits, the network agent must develop an efficient QEC strategy from scratch solely by interacting with the quantum memory at every time step via a set of unitary gates (such as CNOTs and bit-flips) and measurements. They are chosen according to the available hardware and define the action set of the agent. Importantly, the network must react and adapt its strategy to the binary measurement results, providing real-time quantum feedback.

This particular task seems practically unsolvable for the present reinforcement learning techniques if no extra precautions are taken. The basic challenge is also encountered in other difficult RL applications: **the first sequence leading to an increased return is rather long**. In our scenarios, the probability to randomly select a good sequence is much less than  $10^{-12}$ . Moreover, any subsequence may be worse than the trivial (idle) strategy: for example, performing an incomplete encoding sequence (ending up in a fragile entangled state) can accelerate decay. Adopting the straightforward return, namely the overlap of the final and initial states, both the trivial strategy and the error-correction strategy are fixed points. These are separated by a wide barrier – all the intermediate-length sequences with lower return. **In our numerical experiments, naive RL was not successful, except for some tasks with very few qubits and gates.**

We introduce two key concepts to solve this challenge: a two-stage learning approach with one network acting as teacher of another, and a measure of the “recoverable quantum information” retained in any quantum memory.

Before we address these, we mention that from a machine-learning point-of-view there is another unconventional aspect: Instead of sampling initial states of

the RL-environment stochastically, we consider the evolution under the influence of the agent’s actions for all possible states simultaneously. This is required because the quantum memory has to preserve arbitrary input states. Our reward will be based on the completely positive map describing the dissipative quantum evolution of arbitrary states. The only statistical averaging necessary is over measurement outcomes and the probabilistic action choices. Further below, we comment on how this is implemented in practice.

As known from other RL applications (like board games [15]), it helps to provide as much information as possible to the network. In our case, this could mean providing the multi-qubit quantum state at each time step. However, that information is not available in a real experiment. In order to solve this dilemma, we train two different networks in succession (Fig. 2a,b): The first network is fully *state-aware*. Later on, we will use it as a teacher for the second network which essentially only gets the measurement results as an input (plus the information which gate or measurement has been applied). This splits the problem into two sub-problems that are easier to solve. In this approach, the main remaining challenge is to train the state-aware network, while the supervised training of the second network is fairly straightforward in our experience. In contrast, directly training the second network via RL would be tremendously harder, if not impossible, because the input would be significantly less comprehensive than the completely positive map.

At this point, we see that evolving all initial states simultaneously is not only more efficient, but even required to prevent the state-aware network from “cheating”. Otherwise, it might simply memorize the initial state, wait for it to relax, and then reconstruct it – which, of course, is not a valid strategy to preserve a principally unknown quantum state. Such a behavior is avoided when the network is asked to preserve all possible logical qubit states with the same gate sequence. It turns out that this can be implemented efficiently by evolving just four initial quantum states  $\hat{\rho}$  (for a single logical qubit); tracking their evolution fully characterizes, at any point in time, the completely positive map  $\Phi$  of the multi-qubit system that maps  $\hat{\rho}(0)$  to  $\hat{\rho}(t)$ . Moreover, we have found it useful to apply principal component analysis, i. e. to feed only the few largest-weight eigenvectors of the evolved  $\hat{\rho}$ s as input to the network (see Appendix).

We are now ready to define our problem fully from the point of view of reinforcement learning. The state space of the RL environment is the space of completely positive maps. This information is not accessible in a real-world experiment, where the measurements provide partial information about the RL-environment. This reinforcement-learning problem is therefore classified as a partially observed Markov process. This is what is considered in our second learning stage, and our method to solve it relies on a recurrent network. In the modified input scheme of the first learning stage, the agent observes the full state space and we therefore deal with

a fully observed Markov process. In both cases, the RL environment is stochastic due to the measurements. As described above, the action set is defined by the available hardware instructions (unitary gates and measurements).

Two-stage learning with parallel evolution is essential, but not yet sufficient for our challenge. We now introduce a suitable reward that indicates the likely final success of an action sequence ahead of time. In our case, we follow the intuitive idea that this reward should quantify whether the original quantum information survives in the complex entangled many-qubit state that results after application of unitary gates and measurements, and with the system subject to decoherence.

We note that, in the ideal case, without decoherence, two initially orthogonal qubit states are always mapped onto orthogonal states. Therefore, they remain 100% distinguishable, and the original state can always be restored. With a suitable encoding, this remains true even after some errors have happened, if a suitable error-detection and decoding sequence is applied (“recovery”). By contrast, irreversible loss of quantum information means that perfect recovery becomes impossible. In order to make these notions concrete, we start from the well-known fact that the probability to distinguish two quantum states  $\hat{\rho}_1$  and  $\hat{\rho}_2$ , by optimal measurements, is given by the trace distance  $\frac{1}{2}\|\hat{\rho}_1 - \hat{\rho}_2\|_1$ . Let  $\hat{\rho}_{\vec{n}}(t)$  be the quantum state into which the multi-qubit system has evolved, given the initial logical qubit state of Bloch vector  $\vec{n}$ . We now consider the distinguishability of two initially orthogonal states,  $\frac{1}{2}\|\hat{\rho}_{\vec{n}}(t) - \hat{\rho}_{-\vec{n}}(t)\|_1$ . In general, this quantity may display a non-trivial, non-analytic dependence on  $\vec{n}$ . We introduce the “recoverable quantum information” as:

$$\mathcal{R}_Q(t) = \frac{1}{2} \min_{\vec{n}} \|\hat{\rho}_{\vec{n}}(t) - \hat{\rho}_{-\vec{n}}(t)\|_1. \quad (2)$$

The minimum over the full Bloch sphere is taken because the logical qubit state is unknown to the agent, so the success of an action sequence is determined by the worst-case scenario. In other words,  $\mathcal{R}_Q$  specifies a guaranteed value for the remaining distinguishability for all possible logical qubit states. Thus,  $\mathcal{R}_Q$  is a property of the completely positive map that characterizes the dissipative evolution.

The recoverable quantum information  $\mathcal{R}_Q$  is much more powerful than the overlap of initial and final states, as it can be used to construct an immediate reward, evaluating a strategy even at intermediate times. In the idealized case where errors have occurred, but they could *in principle* be perfectly recovered by a suitable detection/decoding sequence,  $\mathcal{R}_Q$  remains 1. As we will see below, this behavior steers the network towards suitable strategies.  $\mathcal{R}_Q$  can be extended towards multiple logical qubits.

As far as  $\mathcal{R}_Q$  is concerned, error correction steps are only required to prevent the multi-qubit system from venturing into regions of the Hilbert space where any further decoherence process would irreversibly destroy the quantum information (and lower  $\mathcal{R}_Q$ ). If one wants the network to actually implement the final decoding sequence, to return back an unentangled state, this can

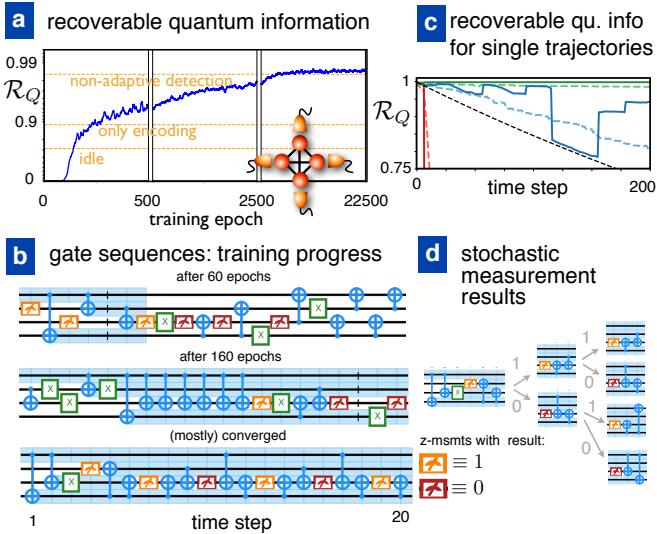


Figure 3. (color) Reinforcement learning for the state-aware network, with a 4-qubit setup (all qubits can be measured, as indicated by the detectors and the red color; all qubit pairs can be subject to CNOT, as indicated by the links). (a) Training progress in terms of the average recoverable quantum information  $\mathcal{R}_Q$ , evaluated at the final step of a 200-step gate sequence. One “epoch” involves training on a batch of 64 trajectories (gate sequences). Eventually, the network performs even better than a combination of encoding and periodic parity checks, due to the adaptive recovery sequences that are triggered by unexpected measurements (i.e. upon error detection). In this example, that leads to ca. 15% increase of the decoherence time over a non-adaptive scheme. (b) Typical gate sequences at different training stages, mostly converged strategy at bottom. Qubits participating in encoding (holding information about the logical qubit state) are indicated with light blue background. (c) Time-evolution of  $\mathcal{R}_Q$ , at the same three different training stages (red, blue, green), for individual trajectories (dashed: averaged over many trajectories). Jumps are due to measurements. (d) Evolution depending on stochastic measurement results, indicated as “0”/“1” and also via the color (red/orange) of the measurement gate symbol. The policy strongly deviates between the different branches, demonstrating that RL finds adaptive quantum feedback strategies, where the behaviour depends on the measurement outcomes. Note that even the mostly converged strategy is still probabilistic to some degree.

be done by adding suitable contributions to the reward (see below).

## Results

We now apply the general approach to different settings, illustrating its flexibility. The training of the state-aware network is analyzed in Fig. 3. In the example, the qubits are subject to bit-flip errors uncorrelated in space and time, with a decay term  $\dot{\hat{\rho}} = T_{\text{dec}}^{-1} \sum_j (\hat{\sigma}_x^{(j)} \hat{\rho} \hat{\sigma}_x^{(j)} - \hat{\rho})$  in the underlying master equation (see Appendix). All of the four qubits may be measured, and there is full

connectivity. During training (Fig. 3a,b), the network first learns to avoid destructive measurements which reveal the logical qubit state. Afterwards, it discovers a gate sequence of CNOTs that creates an entangled state, implementing some version of the 3-qubit repetition code [48, 49]. The particular CNOT sequence shown in the figure generates one possible encoded state out of several equally good ones. The symmetry between these alternative encodings is broken spontaneously during training. The encoding already increases the reward above the trivial level (obtained for storing the logical qubit in one physical qubit only). Finally, the network starts doing repeated parity measurements, of the type  $\text{CNOT}(B \leftrightarrow A)$ ,  $\text{CNOT}(C \leftrightarrow A)$ ,  $M(A)$ , flipping the state of ancilla  $A$  only if the states of  $B$  and  $C$  differ (here  $M$  is a measurement). This implements error detection, helping to preserve the quantum information by preventing the leakage into states with two bit flips that cannot be corrected if undetected. Fig. 3b illustrates the progression from random quantum circuits to a nearly converged strategy. During any single trajectory, the recoverable quantum information can have sudden jumps when measurements are performed (Fig. 3c), with collapses and revivals.

Can we understand better how the network operates? To this end, we visualize the responses of the network responses to the input states (Fig. 4c), projecting the high-dimensional neuron activation patterns into the 2D plane using the t-SNE technique [50]. Similar activation patterns are mapped close to each other, forming clearly visible clusters, each of which results in one type of action. During a gate sequence, the network visits states in different clusters. The sequence becomes complex if unexpected measurement results are encountered (Fig. 4b). In the example shown here, the outcome of the first parity measurement is compatible with three possibilities (one of two qubits has been flipped, or the ancilla state is erroneous). The network has learned to resolve the ambiguity through two further measurements, returning to the usual detection cycle. It is remarkable that RL finds these nontrivial sequences (which would be complicated to construct ab initio), picking out reward differences of a few percent.

The flexibility of the approach is demonstrated by training on different setups, where the network discovers from scratch other feedback strategies (Fig. 5a) adapted to the available resources. For example, we consider a chain of qubits where CNOTs are available only between nearest neighbours and in addition we fix a single measurement location. Then the network learns that it may use the available CNOTs to swap through the chain. However, if every qubit can be measured, the net discovers a better strategy with fewer gates, where the middle two qubits of the chain alternate in playing the role of ancilla. We also show, specifically, the complex recovery sequences triggered by unexpected measurements. They are a-priori unknown, and RL permits to discover them from scratch without extra input. Generally, additional resources (such

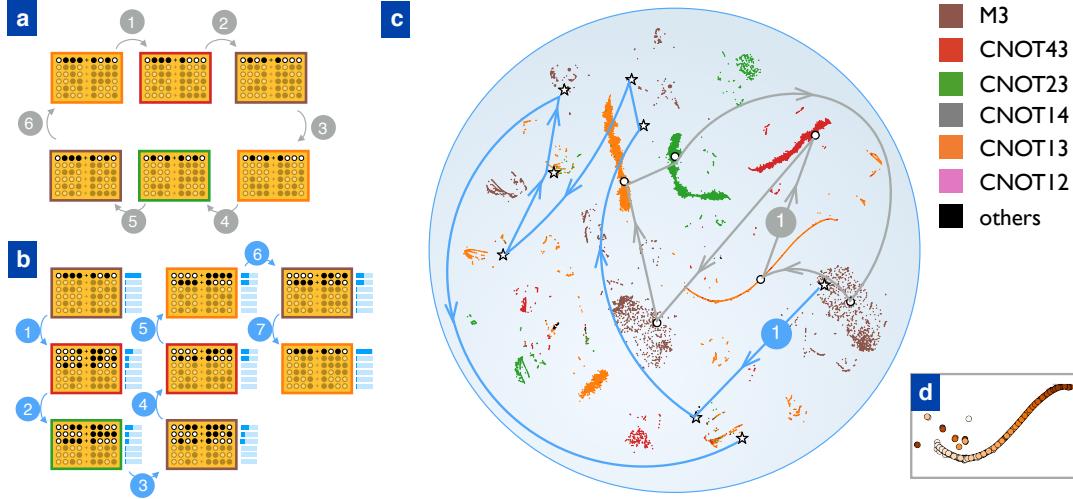


Figure 4. (color) Visualizing the operation of the state-aware network. (Same scenario as in Fig. 3) (a) Sequence of quantum states visited during a standard repetitive detection cycle (after encoding), displayed for an initial logical qubit state  $|+x\rangle = (|0\rangle + |1\rangle)/\sqrt{2}$ , in the absence of unexpected measurement outcomes (no errors). Each state  $\hat{\rho}$  is represented by its decomposition into eigenstates, where e.g.  $\circ \bullet \bullet \bullet + \bullet \circ \bullet \circ \equiv |0111\rangle + |1010\rangle$ . The eigenstates are sorted according to decreasing eigenvalues (probabilities). Eigenstates of less than 5% weight are displayed semi-transparently. In this particular example, each eigenstate is a superposition of two basis states in the z-basis. (b) Gate sequence triggered upon encountering an unexpected measurement. Again, we indicate the states  $\hat{\rho}$ , with bars now showing the probabilities. This sequence tries to disambiguate, by further measurements, the error. (c) Visualization of neuron activations (300 neurons in the last hidden layer), in response to the quantum states (more precisely, maps  $\Phi$ ) encountered in many runs. These activations are projected down to 2D using the t-SNE technique [50], which is a nonlinear mapping that tries to preserve neighborhood relations while reducing the dimensionality. Each of the  $2 \cdot 10^5$  points corresponds to one activation pattern and is colored according to the action taken. The sequences of (a) and (b) are indicated by arrows, with the sequence (b) clearly proceeding outside the dominant clusters (which belong to the more typically encountered states). Qubits are numbered 1, 2, 3, 4, and CNOT13 has control-qubit 1 and target 3. (d) The zoom-in shows a set of states in a single cluster which is revisited periodically during the standard detection cycle; this means we stroboscopically observe the time evolution. The shading indicates the time progressing during the gate sequence, with a slow drift of the state due to decoherence. (cf. Supplementary for extended discussion)

as enhanced connectivity) are exploited to yield better improvement of the decoherence time (Fig. 5b). In another scenario, Fig. 5c, we find that the network successfully learns to adapt to unreliable measurements by redundancy.

In a separate class of scenarios, we consider dephasing of a qubit by a fluctuating field (Fig. 6). If the field is spatially homogeneous and also couples to nearby ancilla qubits, then the dephasing is collective:  $\hat{H}(t) = B(t) \sum_j \mu_j \hat{\sigma}_z^{(j)}$ , where  $B(t)$  is white noise and  $\mu_j$  are the coupling strengths (to qubit and ancillas). Note that, in this situation, one can use neither dynamical decoupling (since the noise is uncorrelated in time) nor decoherence-free subspaces (since the  $\mu_j$  can be arbitrary in general). However, the same RL program used for the examples above also finds solutions here (Fig. 6), without any input specific to the situation (except the available gates). It discovers that the field fluctuations can be tracked and corrected (to some extent) by observing the evolution of the nearby ancillas, measuring them in suitable time intervals. For more than one ancilla, the network discovers a strategy that is adaptive: The choice of measurement

basis depends on the history of previous observations. Brute-force searches in this setting become quickly impossible due to the double-exponential growth of possibilities. The computational effort involved in such a brute-force approach is analyzed in detail in the Supplementary.

Up to now, the network only encodes and keeps track of errors by suitable collective measurements. By revising the reward structure, we can force it to correct errors and finally decode the quantum information back into a single physical qubit. Our objective is to maximize the overlap between the initial and final states, for any logical qubit state (see Appendix). Moreover, we found that learning the decoding during the final time steps is reinforced by punishing states where the logical qubit information is still distributed over multiple physical qubits. The corresponding rewards are added to the previous reward based on the recoverable quantum information. The network now indeed learns to decode properly (Fig. 7a). In addition, it corrects errors. It does so typically soon after detecting an error, instead of at the end of the gate sequence. We conjecture this is because it tries to return as soon as possible back to the known, familiar encoded state. For the same reason, error correction sometimes

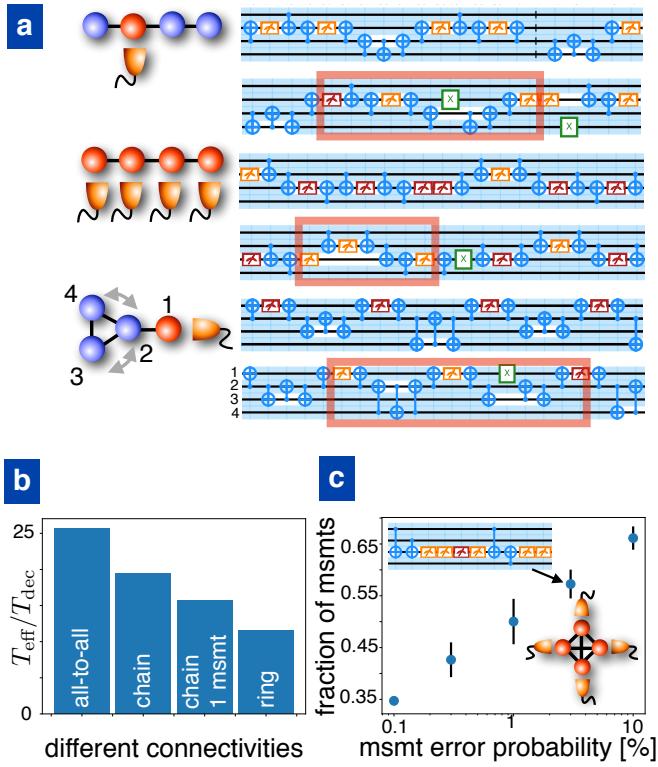


Figure 5. (color) (a) Scenarios with different qubit connectivity. CNOTs are allowed only between qubits connected by a line. In each case, we display the “standard” gate sequence discovered by the net, as well as a sequence involving corrective actions (triggered by an unexpected measurement). From top to bottom: chain with fixed measurement location, chain with arbitrary measurements, ring connected to an ancilla. The red rectangle highlights an interval during which the quantum state is not dominated by a single component, indicating that the precise location of the error still has to be pinpointed. These nontrivial detection/recovery sequences are considerably more complex than the periodic detection cycle. (b) The effective enhancement of the decoherence time via error correction, for the different scenarios. Here,  $T_{\text{dec}} = 1200$  is the single-qubit decoherence time (in units of the gate time that defines the time step), and  $T_{\text{eff}}$  was extracted from the decay of  $\mathcal{R}_Q$  after 200 time steps. The differences can be traced back to the lengths of the detection cycles. (c) Behaviour as a function of measurement error. The network discovers that redundancy is needed, i. e. the number of measurements in the gate sequences increases (in this plot, from 1 per cycle to about 6).

even happens without an explicit reward.

So far, we have trained the state-aware network. However, this cannot yet be applied to an experiment, where the quantum state is inaccessible to us. This requires a network whose only input consists in the measurement results (and the selected gates, since the policy is probabilistic), requiring some sort of memory. An elegant solution consists in a *recurrent* neural network. We use the widespread long short-term memory (LSTM) approach [51].

Once the first, state-aware network has been trained

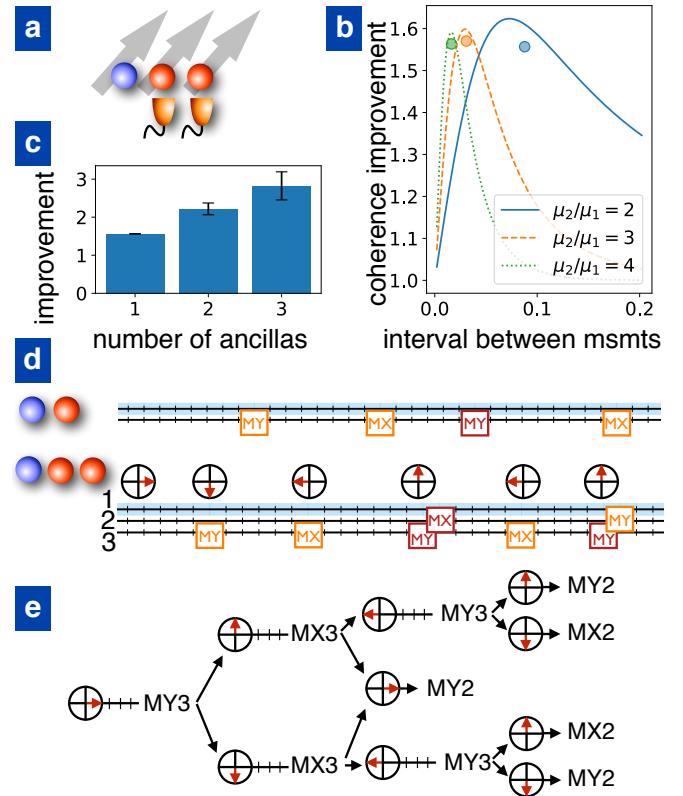


Figure 6. (color) Counteracting dephasing by measurements (adaptive phase estimation). (a) The setting: a data qubit, whose phase undergoes a random walk due to a fluctuating field. Since the field is spatially correlated, its fluctuations can be detected by measuring the evolution of nearby ancilla qubits, which can then be exploited for correction. In this example, the allowed actions are measurements along x,y (MX, MY), and the idle operation. (b) For one ancilla, the network performs measurements periodically, finding the optimal measurement interval. This can be seen by comparing the coherence time enhancement as a function of the interval (in units of the single-qubit decoherence time  $T_{\text{single}}$ ) for the network (circles) with the analytical predictions (curves; cf. Supplementary). Remaining differences are due to the discretization of the interval (not considered in the analytics). The coupling between noise and ancilla ( $\mu_2$ ) is different from that between noise and data qubit ( $\mu_1$ ), and the strategy depends on the ratio indicated here. (c) Coherence time enhancement, for different numbers of ancillas (here  $\mu_2 = \mu_3 = \mu_4 = 4\mu_1$ ). (d) Gate sequences. For two ancillas, the network discovers an adaptive strategy, where measurements on qubit #2 are rare, and the measurement basis is decided based on previous measurements of qubit #3. The arrows show the measurement result for qubit #3, in the equator plane of the Bloch sphere. (e) The 2-ancilla adaptive strategy (overview; see also Supplementary). Here, a brute-force search for strategies is still (barely) possible, becoming infeasible for higher ancilla numbers due to the exponentially large search space of adaptive strategies. [ $\mu_2 = 3.8\mu_1$ ,  $\mu_3 = 4.1\mu_1$  in (d),(e)]

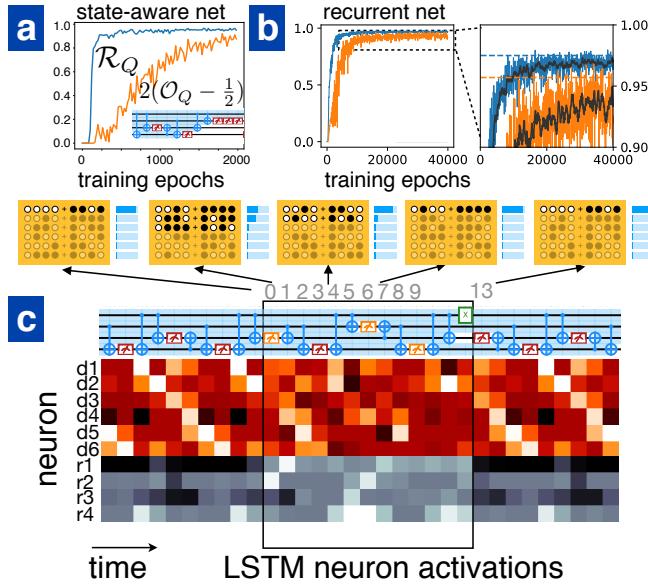


Figure 7. (color) (a) Training the state-aware network to implement decoding back into the original data qubit, near the end of a 200-step gate sequence. Blue: recoverable quantum information  $\mathcal{R}_Q$ . Orange: (rescaled, shifted) overlap  $\mathcal{O}_Q$  between final and initial state of the data qubit. Inset displays the decoding gate sequence. (b) Training the recurrent network in a supervised way to imitate the state-aware network. Again, we show  $\mathcal{R}_Q$  and  $\mathcal{O}_Q$  evolving during training. Dashed blue/orange lines depict the performance of the state-aware network. (c) To investigate the workings of the recurrent network, we display some of the LSTM neuron activations in the second-to-last layer. Quantum states are illustrated even though the network is unaware of them. The network's input consists of the observed measurement result (if any) and of the actual gate choice (made in the previous time step, here aligned on top of the current neuron activations). The example gate sequence displayed here first shows the repetitive standard error detection pattern of repeated parity measurements that is interrupted once an unexpected measurement is encountered, indicating an error. During error recovery (boxed region), the network first tries to pinpoint the error and then applies corrections. The neurons whose behavior changes markedly during recovery are depicted in gray. Neuron 1r obviously keeps track of whether recovery is ongoing, while 3r acts like a counter keeping time during the standard periodic detection pattern. (see Supplementary for more)

successfully, it is used as a teacher in supervised learning to train the second, recurrent network (Fig. 7b). This could then be applied as a controller to experimental runs, deciding on gate sequences depending on measurements. It might also be refined by RL, e.g. to adapt to changes in the parameters (decoherence rates etc.). Learning to correct (see above) is essential for successful training of the recurrent network, since the latter must learn to consider measurement results distributed in time and deduce the proper corrective actions.

We have trained the recurrent network based on a fully converged state-aware network. Inspecting the LSTM

neuron activations (Fig. 7c), we see that different neurons activate for different events and some clearly display prolonged memory (remaining active during certain time-intervals relevant for the strategy). For example, one neuron switches on during the recovery sequence after an unexpected measurement, while another seems like an internal counter operating during the periodic detection sequence.

We now come back to the statement in the introduction that our approach is fully autonomous and can be applied to a broad range of problems with small human effort. In all the preceding examples, and also in general, the only human input to our approach is the problem specification, primarily the noise model (specifying the dissipative time evolution governing the quantum state) and the particular action set (i.e., the available hardware instructions related to the setup and its connectivity). Importantly, fine-tuning the hyperparameters (like learning rate, network architecture, etc.) is not required; in the Supplementary, we demonstrate that a common set of hyperparameters can be used for all the scenarios.

### Possible future applications

The physical setups considered in today's quantum computing platforms contain many components and features that go beyond the simplest scenario of short-range coupled qubits. Conceptually, the approach developed in the present work is general enough to find future application in any of the following experimentally relevant domains.

An important example is cavities, which can be used as long-lived quantum memory, especially in the microwave domain. When they are coupled to qubits, nonlinear operations can be performed that may aid in error correction of the cavity state, as the Yale group has demonstrated ("kitten" and "cat" codes [52, 53]). Our approach allows to cover such situations without any changes to the reinforcement learning method. Only the description of the physical scenario, via the set of available actions, and of course the physics simulation will have to be updated. Cavities also give access to unconventional controls, e.g. naturally occurring long-distance multi-qubit entangling gates provided by the common coupling of the qubits to the cavity. In addition, they permit direct collective readout that is sensitive to the joint state of multiple qubits, which may be used to speed up error detection operations. Again, RL based quantum feedback of the type proposed here can naturally make use of these ingredients.

Novel hardware setups, like cross-bar type geometries [54, 55], give rise to the challenge to exploit the unconventional connectivity, for which our approach is well suited. In the future, it may even become possible to co-optimize the hardware layout (taking into account physical constraints) and the strategies adapted to the layout. In the simplest case, this means discovering strategies for automatically generated alternative layouts and comparing their performance.

The actions considered by the agent need not refer to unitary operations. They might also perform other functions, like restructuring the connectivity itself in real-time. This is the case for the proposed 2D ion-trap architecture where the ions are shuffled around using electrodes [19]. Similar ideas have been proposed for spins in quantum dots, which can be moved around using electrodes or surface-acoustic waves. Again, no changes to our approach would be needed. The modifications are confined to the physics simulation. Depending on the present state of the connectivity, the set of effective qubit gates would change.

Like any numerical approach, our method is invariably limited to modest qubit numbers (of course, these will increase with further optimizations, possibly up to about 10). It is important, therefore, to recall that even an improvement of the decoherence rate in an isolated few-qubit module can have useful applications (as a quantum memory, e.g. in a quantum repeater). More generally, it is clear that classical simulation of a full-scale quantum computer in the domain of quantum supremacy is out of the question, by definition. This is a challenge widely acknowledged by the entire community, affecting not only optimization but also design, testing, and verification of a quantum computer. One promising way to address this challenge at least partially, advocated by a growing number of experimental groups, is the so-called modular approach to quantum computation and quantum devices. This consists in connecting small few-qubit quantum modules together via quantum network links [19, 20]. The main advantage of this approach is the ability to control and debug small quantum modules as opposed to an entire large monolithic quantum computer. Our approach is very well suited to this strategy. In principle one can even envision a hierarchical application of the quantum module concept (with error correction strategies applied to multiple modules coupled together), but for that case our approach would need to be extended (e.g. by using RL to find one- and two-qubit gates acting on the logical qubits stored inside the modules).

## Conclusions

We have seen how a network can discover quantum error correction techniques from scratch. It finds a-priori unknown nontrivial detection/recovery sequences for diverse settings without any more input than the available gate set. The trained neural networks can in principle be used to control experimental quantum devices. The present approach is flexible enough to be applied directly to a range of further, qualitatively different physical situations, like non-Markovian noise, weak measurements, qubit-cavity systems, and error-corrected transport of quantum information through networks. An obvious challenge for the future is to successfully discover strategies on even more qubits, where eventually full protection against all noise sources and multiple logical qubits could be re-

alized. There is still considerable leeway in improving the speed of the physics simulation and of GPU-based training (for further details on the current computational effort, see appendix L).

On the machine learning side, other RL schemes can be substituted for the natural policy gradient adopted here, like Q-learning or advantage-actor-critic techniques, or RL with continuous controls. Recurrent networks might be employed to discover useful subsequences. The two-stage learning approach introduced here could also be applied in other RL scenarios, where one would first train based on expanded state information. In general, we have shown that neural-network based RL promises to be a flexible and general tool of wide-ranging applicability for exploring feedback-based control of quantum and classical systems in physics.

**Acknowledgements** We thank Hugo Ribeiro and Vittorio Peano for fruitful comments on the manuscript.

**Author Contributions** All authors contributed to the ideas, their implementation, and the writing of the manuscript. The numerics was performed by T.F., P.T. and T.W.

**Data availability** The data that support the plots within this paper and other findings of this study are available from the corresponding author on request.

*Note added:* Shortly before submission of the present manuscript to the arXiv, a preprint [56] appeared exploring RL with recurrent networks for optimal quantum control (without feedback).

## A. Physical time evolution

To track the time evolution for an *arbitrary* initial logical qubit state (identified by its Bloch vector  $\vec{n}$ ), we start from  $\hat{\rho}_{\vec{n}}(0) = \frac{1}{2}(\mathbb{1} + \vec{n}\hat{\sigma}) \otimes \hat{\rho}_{\text{Rest}}$ , factorizing out the (fixed) state of all the other qubits. Now consider the four quantities

$$\hat{\rho}_0(0) = \frac{1}{2}(\hat{\rho}_{\vec{e}_j}(0) + \hat{\rho}_{-\vec{e}_j}(0)) \quad (\text{A1a})$$

$$\delta\hat{\rho}_j(0) = \frac{1}{2}(\hat{\rho}_{\vec{e}_j}(0) - \hat{\rho}_{-\vec{e}_j}(0)) \quad (\text{A1b})$$

where  $j \in \{x, y, z\}$  and  $\vec{e}_j$  are the basis vectors; note that the right-hand side of Eq. (A1a) is independent of  $j$ .  $\hat{\rho}_0$  and the  $\delta\hat{\rho}_j$  are evolved stepwise, for each time-interval  $[t_i, t_f]$  according to the update rule

$$\hat{\rho}_0(t_f) = \frac{\phi[\hat{\rho}_0(t_i)]}{\text{tr}(\phi[\hat{\rho}_0(t_i)])} \quad (\text{A2a})$$

$$\delta\hat{\rho}_j(t_f) = \frac{\phi[\delta\hat{\rho}_j(t_i)]}{\text{tr}(\phi[\hat{\rho}_0(t_i)])}. \quad (\text{A2b})$$

In the absence of measurements,  $\phi$  is the completely positive map for the given time-interval. In the presence of measurements, it is an unnormalized version (see below). We explicitly renormalize such that always  $\text{tr}(\hat{\rho}_0(t)) = 1$ .

$\hat{\rho}_0$  and the  $\delta\hat{\rho}_j$  give us access to the density matrix for every logical qubit state, at any time  $t$ :

$$\hat{\rho}_{\vec{n}}(t) = \frac{\hat{\rho}_0(t) + \sum_j n_j \delta\hat{\rho}_j(t)}{1 + \text{tr}(\sum_j n_j \delta\hat{\rho}_j(t))} \quad (\text{A3})$$

## B. Physical scenarios

We always start from the initial condition that the logical qubit is stored in one physical qubit, and the others are prepared in the down state ( $|1\rangle$ ). If explicit recovery is desired, we use this original qubit also as the target qubit for final decoding. The time evolution is divided into discrete time steps of uniform length  $\Delta t$  (set to 1 in the main text). At the start of each of these time slices, we perform the measurement or gate operation (which is assumed to be quasi-instantaneous) chosen by the agent; afterwards, the system is subject to the dissipative dynamics. Thus, the map  $\phi$  for the time interval  $[t, t + \Delta t]$  is of the form  $\phi[\hat{\rho}] = e^{\Delta t \mathcal{D}}(\hat{U}\hat{\rho}\hat{U}^\dagger)$  for unitary operations  $\hat{U}$  and  $\phi[\hat{\rho}] = e^{\Delta t \mathcal{D}}(\hat{P}_m\hat{\rho}\hat{P}_m^\dagger)$  for projection operators  $\hat{P}_m$  ( $m$  indicates the measurement results) where  $\mathcal{D}$  is the dissipative part of the Liouvillian (we only consider Markovian noise). Note that the measurement results are chosen stochastically according to their respective probability  $\text{tr}(\hat{P}_m\hat{\rho}_0)$ . In the examples discussed in the figures, we use two different error models, the bit-flip error (bf) and the correlated noise error (cn):

$$\mathcal{D}_{\text{bf}}\hat{\rho} = T_{\text{dec}}^{-1} \sum_q \hat{\sigma}_x^{(q)} \hat{\rho} \hat{\sigma}_x^{(q)} - \hat{\rho} \quad (\text{B1a})$$

$$\mathcal{D}_{\text{cn}}\hat{\rho} = T_{\text{dec}}^{-1} \left( \hat{L}_{\text{cn}}\hat{\rho}\hat{L}_{\text{cn}}^\dagger - \frac{1}{2} \left\{ \hat{L}_{\text{cn}}^\dagger \hat{L}_{\text{cn}}, \hat{\rho} \right\} \right) \quad (\text{B1b})$$

where  $\hat{\sigma}_{x,y,z}^{(q)}$  applies the corresponding Pauli operator to the  $q$ th qubit and  $\hat{L}_{\text{cn}} = \frac{1}{\sqrt{\sum_q \mu_q^2}} \sum_q \mu_q \hat{\sigma}_z^{(q)}$ . Here,  $\mu_q$  denotes the coupling of qubit  $q$  to the noise. Note that in the bit-flip scenario the single qubit decay time  $T_{\text{single}} = T_{\text{dec}}$ , whereas in the presence of correlated noise it is  $T_{\text{single}} = T_{\text{dec}} \left( \sum_q \mu_q^2 \right) / \mu_1^2$ .

## C. Recoverable quantum information

Based on Eq. (A3),  $\mathcal{R}_Q$  as introduced in the main text can be written as

$$\mathcal{R}_Q(t) = \min_{\vec{n}} \left\| \sum_j n_j \delta\hat{\rho}_j(t) \right\|_1 \quad (\text{C1})$$

in the (for us relevant) case that  $\text{tr}(\delta\hat{\rho}_j(t)) = 0$  for all  $j$ .

The trace distance  $\frac{1}{2} \left\| \sum_j n_j \delta\hat{\rho}_j(t) \right\|_1$  has often a non-trivial dependence on the logical qubit state  $\vec{n}$  and finding its minimum can become nontrivial. However, the location of the minimum can sometimes be “guessed”

in advance. For any CHZ quantum circuit [48], i.e. for all the bit-flip examples considered here, the anti-commutator relation  $\{\delta\hat{\rho}_j, \delta\hat{\rho}_k\} = 0$  is satisfied for all distinct  $j \neq k$ ; it can be shown that this restricts the minimum to lie along one of the coordinate axes:  $\min_{\vec{n}} \left\| \sum_j n_j \delta\hat{\rho}_j(t) \right\|_1 = \min_{j \in \{x,y,z\}} \|\delta\hat{\rho}_j(t)\|_1$ . For the correlated noise, the trace distance  $\frac{1}{2} \left\| \sum_j n_j \delta\hat{\rho}_j(t) \right\|_1$  is symmetric around the  $z$ -axis and takes its minimal value at the equator.

After a measurement, the updated value of  $\mathcal{R}_Q$  may vary between the different measurement results. To obtain a measure that does not depend on this, we introduce  $\bar{\mathcal{R}}_Q$  as the average over all possible values of  $\mathcal{R}_Q$  (after a single time step), weighted by the probability to end up in the corresponding branch. If the action is not a measurement, there is only one option and thus  $\bar{\mathcal{R}}_Q = \mathcal{R}_Q$ .

## D. Protection reward

The goal of the “protection reward” is to maximize  $\mathcal{R}_Q$  at the end of the simulation, i.e. the ability to *in principle* recover the target state. A suitable (immediate) reward is given by

$$r_t = \begin{cases} 1 + \frac{\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t)}{2\Delta t / T_{\text{single}}} + 0 & \text{if } \bar{\mathcal{R}}_Q(t+1) > 0 \\ 0 & \text{if } \mathcal{R}_Q(t) \neq 0 \\ 0 & \underbrace{-P}_{\text{if } \bar{\mathcal{R}}_Q(t+1) = 0} \\ 0 & \underbrace{+0}_{\text{if } \mathcal{R}_Q(t) = 0} \end{cases} \quad (\text{D1})$$

with  $\mathcal{R}_Q$  and  $\bar{\mathcal{R}}_Q$  as defined above,  $T_{\text{single}}$  the decay time for encoding in one physical qubit only (“trivial” encoding),  $\Delta t$  the time step, and  $P$  a punishment for measurements which reveal the logical qubit state. Based on this reward, we choose the return (the function of the reward sequence used to compute the policy gradient) as

$$R_t = (1 - \gamma) \left( \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}^{(1)} \right) + r_t^{(2)} \quad (\text{D2})$$

where  $\gamma$  is the return discount rate; for more information on the (discounted) return, see e.g. [57].

## E. Recovery reward

The protection reward does not encourage the network to finally decode the quantum state. If this behavior is desired, we add suitable terms to the reward (only employed for Fig. 7):

$$r_t^{(\text{recov})} = \beta_{\text{dec}} (D_{t+1} - D_t) + \beta_{\text{corr}} C_T \delta_{t,T-1} \quad (\text{E1})$$

where  $D_t = \frac{1}{2} (I_t^{(q')} - \sum_{q \neq q'} I_t^{(q)})$ , unless  $t \leq T_{\text{signal}}$  where we set  $D_t = 0$ . This means decoding is only rewarded after  $T_{\text{signal}}$ . We set  $I_t^{(q)} = 1$  if  $\text{tr}_{\vec{q}}(\delta\hat{\rho}_j(t)) \neq 0$  for any  $j$ ,

and otherwise  $I_t^{(q)} = -1$ .  $\text{tr}_{\bar{q}}$  denotes the partial trace over all qubits except  $q$ , and  $q'$  labels the target qubit. The condition  $I_t^{(q)} = 1$  implies that the logical qubit state is encoded in the specific qubit  $q$  (this is not a necessary criterion).  $C_T$  is 1 if (at the final time  $T$ ) the logical qubit state is encoded in the target qubit only and this qubit has the prescribed polarization (i.e. not flipped), and otherwise 0.

As return, we use

$$R_t^{(\text{recov})} = (1 - \gamma) \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}^{(\text{recov})} \quad (\text{E2})$$

with the same return discount rate  $\gamma$  as for the protection reward.

With this reward, we aim to optimize the minimum overlap  $\mathcal{O}_Q = \min_{\vec{n}} \langle \phi_{\vec{n}} | \text{tr}_{\bar{q}'}(\hat{\rho}_{\vec{n}}) | \phi_{\vec{n}} \rangle$  between the (pure) target state  $|\phi\rangle_{\vec{n}}$  and the actual final state  $\hat{\rho}_{\vec{n}}$  reduced to the target qubit, given by the partial trace  $\text{tr}_{\bar{q}'}(\hat{\rho}_{\vec{n}})$  over all other qubits.

## F. Input of the state-aware network

The core of the input to the state-aware network is a representation of the density matrices  $\hat{\rho}_0$ ,  $\hat{\rho}_1 := \hat{\rho}_0 + \delta\hat{\rho}_x$ ,  $\hat{\rho}_2 := \hat{\rho}_0 + \delta\hat{\rho}_y$ , and  $\hat{\rho}_3 := \hat{\rho}_0 + \delta\hat{\rho}_z$ . Together, they represent the completely positive map of the evolution (for arbitrary logical qubit states). For reduction of the input size (especially in view of higher qubit numbers), we compress them via principal component analysis (PCA), i.e. we perform an eigendecomposition  $\hat{\rho}_j = \sum_k p_k |\phi_k\rangle\langle\phi_k|$  and select the eigenstates  $|\phi_k\rangle$  with the largest eigenvalues  $p_k$ . To include also the eigenvalue in the input, we feed all components of the scaled states  $|\tilde{\phi}_k\rangle := \sqrt{p_k} |\phi_k\rangle$  (which yield  $\hat{\rho}_j = \sum_k |\tilde{\phi}_k\rangle\langle\tilde{\phi}_k|$ ) into the network, where the states are in addition sorted by their eigenvalue. For our simulations, we select the 6 largest components, so we need  $768 = 4 \cdot 6 \cdot 16 \cdot 2$  input neurons (4 density matrices, 16 is the dimension of the Hilbert space, 2 for real and imaginary part).

In addition, at each time step we indicate to the network whether a potential measurement would destroy the quantum state by revealing the quantum information. Explicitly, we compute for each measurement whether  $\text{tr}(\hat{P}\delta\hat{\rho}_j) = 0$  for all  $j \in \{x, y, z\}$  and every possible projector  $\hat{P}$  (i.e. every possible measurement result), and feed these boolean values into the network. Note that this information can be deduced from the density matrix (so in principle the network could learn that deduction on its own, but giving it directly speeds up training).

Because all relevant information for the decision about the next action is contained in the current density matrix, knowledge about the previous actions is not needed. However, we have found that providing this extra information is helpful to accelerate learning. Therefore, we provide also the last action (in a one-hot encoding). We note

that it is not necessary to feed the latest measurement result to the network, since the updated density matrix is conditional on the measurement outcome and therefore contains all relevant information for future decision.

To train the state-aware network to restore the original state at the end of a trajectory, it becomes necessary to add the time to the input. It is fully sufficient to indicate the last few time steps where  $t > T_{\text{signal}}$  (when decoding should be performed) in a one-hot encoding.

## G. Layout of the state-aware network

Our state-aware networks have a feedforward architecture. Between the input layer and the output layer (one neuron per action), there are two or three hidden layers (the specific numbers are summarized in the last section of the Appendix). All neighboring layers are densely connected, the activation function is the rectified linear unit (ReLU). At the output layer, the softmax function  $x_j \mapsto e^{x_j} / \sum_k e^{x_k}$  is applied such that the result can be interpreted as a probability distribution.

## H. Reinforcement learning of the state-aware network

Our learning scheme is based on the policy gradient algorithm [45]. The full expression for our learning gradient (indicating the change in  $\theta$ ) reads

$$g = \lambda_{\text{pol}} F^{-1} \cdot \mathbb{E} \left[ \sum_t (R_t - b_t) \frac{\partial}{\partial \theta} \ln \pi_\theta(a_t | s_t) \right] - \lambda_{\text{entr}} \mathbb{E} \left[ \sum_a \frac{\partial}{\partial \theta} [\pi_\theta(a | s) \ln \pi_\theta(a | s)] \right] \quad (\text{H1})$$

where  $R_t$  is the (discounted) return (cmp. Eqs. (D2) and (E2)). This return is corrected by an (explicitly time-dependent) baseline  $b_t$  which we choose as exponentially decaying average of  $R_t$ , i.e. for the training update in epoch  $N$ , we use  $b_t = (1 - \kappa) \sum_{n=0}^{N-1} \kappa^n \bar{R}_t^{(N-1-n)}$  where  $\kappa$  is the baseline discount rate and  $\bar{R}_t^{(n)}$  is the mean return at time step  $t$  in epoch  $n$ . We compute the natural gradient [58–60] by multiplying  $F^{-1}$ , the (Moore–Penrose) inverse of the Fisher information matrix  $F = \mathbb{E}[(\frac{\partial}{\partial \theta} \ln \pi_\theta(a | s))(\frac{\partial}{\partial \theta} \ln \pi_\theta(a | s))^T]$ . The second term is entropy regularization [61]; we use it only to train the state-aware network shown in Fig. 7. As update rule, we use adaptive moment estimation (Adam [62]) without bias correction.

## I. Layout of the recurrent network

The recurrent network is designed such that it can in principle operate in a real-world experiment. This means

in particular that (in contrast to the state-aware network) its input must not contain directly the quantum state (or the evolution map); instead, measurements are its only way to obtain information about the quantum system. Hence, the input to the recurrent network contains the present measurement result (and additionally the previous action). Explicitly, we choose the input as a one-hot encoding for the action in the last time step, and in case of measurements, we additionally distinguish between the different results. In addition, there is an extra input neuron to indicate the beginning of time (where no previous action was performed). Since this input contains only the most recent “event”, the network requires a memory to perform reasonable strategies, i.e. we need a recurrent network. Therefore, the input and output layer are connected by two successive long short-term memory (LSTM) layers [51] with tanh inter-layer activations. After the output layer, the softmax function is applied (like for the state-aware network).

### J. Supervised learning of the recurrent network

The training data is generated from inference of a state-aware network which has been trained to sufficiently good strategies (via reinforcement learning); for every time step in each trajectory, we save the network input and the policy, i.e. the probabilities for all the actions, and we train on this data. It is possible to generate enough data such that overfitting is not a concern (for the example in Fig. 7, each trajectory is reused only 5 times during the full training process). For the actual training of the recurrent network, we use supervised learning with categorical cross-entropy as cost function ( $q$  is the actual policy of the recurrent network to train, and  $p$  the desired policy from the state-aware network):

$$C(q, p) = - \sum_a p(a) \ln q(a) \quad (\text{J1})$$

Due to the LSTM layers, it is necessary to train on full trajectories (in the true time sequence) instead of individual actions. Dropout [63] is used for regularization. The training update rule is adaptive moment estimation (Adam [62]).

### K. Physical parameters and hyperparameters

The physical parameters used throughout the main text are summarized in the following table. Times are always given in units of the time step (gate time).

#### Physical parameters

Figs. 3, 4, 5, 7 decoherence time $T_{\text{dec}}$	1200
Fig. 6 single qubit decoherence time	500
$T_{\text{single}} = T_{\text{dec}}/\mu_1^2$	
Figs. 3, 4, 5, 7 number of time steps $T$	200
Fig. 6 number of time steps $T$	100

We have used a few separately trained neural network throughout this work which differ slightly in hyperparameters (e.g. in the number of hidden layers and neurons per layer). This is not due to fine-tuning, and in the Supplementary Information we demonstrate that we can successfully train the neural networks in all scenarios with one common set of hyperparameters. The strategies found by the neural networks are not influenced by using different sets of hyperparameters. Different hyperparameters may influence the training time etc. In the following table, we summarize the architecture of the networks, i.e., we list the number of neurons in each layer.

Network architectures	Neurons per layer
Fig. 3	(793, 300, 300, 21)
Fig. 4	(793, 300, 300, 300, 21)
Fig. 5a-Fig. 5b all-to-all connected:	(793, 300, 300, 21)
Fig. 5a-Fig. 5b chain:	(787, 300, 300, 300, 15)
Fig. 5a-Fig. 5b chain with one msmt qubit:	(781, 600, 300, 12)
Fig. 5a-Fig. 5b circle:	(783, 300, 300, 300, 14)
Fig. 5c (all identical)	(793, 300, 300, 300, 21)
Fig. 6b-Fig. 6d: one ancilla	(101, 100, 50, 3)
Fig. 6c, Fig. 6d: two ancillas	(265, 100, 50, 5)
Fig. 6c: three ancillas	(781, 200, 100, 7)
Fig. 7a	(803, 300, 300, 21)
Fig. 7b-Fig. 7c (hidden layers with LSTM units)	(26, 128, 128, 21)

Each output neuron represents one action that can be performed by the agent and thus, the output layer size is equal to the number of actions. In the bit-flip scenarios (Figs. 3, 4, 5, 7), the actions are CNOTs according to connectivity, measurements along  $z$  as indicated in the corresponding sketches, deterministic bit flips on each qubit, and idle. When dealing with correlated noise, cf. Fig. 6, the available actions are instead measurements along  $x$  and  $y$  on all ancilla qubits, and the idle operation. Note that our whole approach is general and able in principle to deal with arbitrary quantum gates.

The hyperparameters used for training the state-aware networks are summarized in the following table:

**Hyperparameters  
state-aware network**

training batch size	64	
Adam parameters $\eta, \beta_1, \beta_2$	0.0003 <sup>a</sup> , 0.9, 0.999 (no bias correction)	
return discount rate $\gamma$	0.95	
baseline discount rate $\kappa$	0.9	
punishment reward coefficient $P$	0.1	
decoding reward coefficient	20	
$\beta_{\text{dec}}$ (only used in Fig. 7a)		
correction reward coefficient	10	
$\beta_{\text{corr}}$ (only used in Fig. 7a)		
decoding signal time $T_{\text{signal}}$ (only used in Fig. 7a)	$T_{\text{signal}} = T - 10 = 190$	
reward scale $\lambda_{\text{pol}}$	4.0	
entropy regularization $\lambda_{\text{entr}}$ (only used in Fig. 7a)	$\lambda_{\text{entr}} = 5 \cdot 10^{-3}$ for the first 12 000 training epoches, $\lambda_{\text{entr}} = 0$ afterwards	

<sup>a</sup> The exact value for the learning rate used in the simulations is in fact  $0.0001\sqrt{10}$ ; the irrational factor of  $\sqrt{10}$  is caused by a slight deviation between our implementation and the standard Adam scheme which in the end resulted only in a redefinition of the learning rate.

The hyperparameters used for training the recurrent network (cf. Fig. 7b and Fig. 7c) are summarized in the following table:

**Hyperparameters  
recurrent network**

training batch size	16
Adam parameters $\eta, \beta_1, \beta_2$	0.001, 0.9, 0.999
dropout level (after each LSTM layer)	0.5

Our RL implementation relies on the Theano framework [64] (and Keras for defining networks).

**L. Computational resources**

The computationally most expensive tasks in this paper are the training runs of the state-aware networks depicted in Figs. 3, 5 and 7a. Full training for a given scenario can be achieved using 1.6 million training trajectories, which can be run within 6 hours on a single CPU+GPU node (CPU: Intel Xeon E5-1630 v4, GPU: Nvidia Quadro P5000). Currently, more than 2/3 of the time are spent on the numerical simulation of the physics time evolution, which is still performed on the CPU. We expect that the total runtime can be improved significantly by a more efficient implementation and more powerful hardware (including also an implementation of the physics simulation on GPU). The memory consumption is very modest (for these examples below 200 MB), dominated mostly by the need for storing the input to the network for all trajectories inside a batch (here, 64 trajectories with 200 time steps each) and less by the network weights (here, up to about 600,000). For a more detailed discussion, we refer the reader to the supplementary material (sec. 6).

- 
- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015.
  - [2] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. CreateSpace Independent Publishing Platform, September 2016. Google-Books-ID: PQI7vgAACAAJ.
  - [3] Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325):602–606, February 2017.
  - [4] Juan Carrasquilla and Roger G. Melko. Machine learning phases of matter. *Nat Phys*, 13(5):431–434, May 2017.
  - [5] Evert P. L. van Nieuwenburg, Ye-Hua Liu, and Sebastian D. Huber. Learning phase transitions by confusion. *Nature Physics*, 13(5):435–439, May 2017.
  - [6] Giacomo Torlai and Roger G. Melko. Neural Decoder for Topological Codes. *Physical Review Letters*, 119(3):030501, July 2017.
  - [7] Moritz August and Xiaotong Ni. Using Recurrent Neural Networks to Optimize Dynamical Decoupling for Quantum Memory. *Physical Review A*, 95(1), January 2017. arXiv: 1604.00279.
  - [8] Vedran Dunjko and Hans J. Briegel. Machine learning and artificial intelligence in the quantum domain. *arXiv:1709.02779 [quant-ph]*, September 2017. arXiv: 1709.02779.
  - [9] P. Baireuther, T. E. O’Brien, B. Tarasinski, and C. W. J. Beenakker. Machine-learning-assisted correction of correlated qubit errors in a topological code. *arXiv:1705.07855 [cond-mat, physics:quant-ph]*, May 2017. arXiv: 1705.07855.
  - [10] Stefan Krastanov and Liang Jiang. Deep Neural Network Probabilistic Decoder for Stabilizer Codes. *Scientific Reports*, 7(1):11003, September 2017.
  - [11] Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre G. R. Day, Clint Richardson, Charles K. Fisher, and David J. Schwab. A high-bias, low-variance introduction to Machine Learning for physicists. *arXiv:1803.08823 [cond-mat, physics:physics, stat]*, March 2018. arXiv: 1803.08823.
  - [12] Ross D. King, Jem Rowland, Stephen G. Oliver, Michael Young, Wayne Aubrey, Emma Byrne, Maria Liakata, Magdalena Markham, Pinar Pir, Larisa N. Soldatova, Andrew Sparkes, Kenneth E. Whelan, and Amanda Clare. The Automation of Science. *Science*, 324(5923):85–89, April 2009.
  - [13] Michael Schmidt and Hod Lipson. Distilling Free-Form Natural Laws from Experimental Data. *Science*, 324(5923):81–85, April 2009.
  - [14] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex

- Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.
- [15] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354, October 2017.
- [16] C. Chen, D. Dong, H. X. Li, J. Chu, and T. J. Tarn. Fidelity-Based Probabilistic Q-Learning for Control of Quantum Systems. *IEEE Transactions on Neural Networks and Learning Systems*, 25(5):920–933, May 2014.
- [17] Marin Bukov, Alexandre G. R. Day, Dries Sels, Phillip Weinberg, Anatoli Polkovnikov, and Pankaj Mehta. Machine Learning Meets Quantum State Preparation. The Phase Diagram of Quantum Control. *arXiv:1705.00565 [cond-mat, physics:quant-ph]*, May 2017. arXiv: 1705.00565.
- [18] Alexey A. Melnikov, Hendrik Poulsen Nautrup, Mario Krenn, Vedran Dunjko, Markus Tiersch, Anton Zeilinger, and Hans J. Briegel. Active learning machine learns to create new quantum experiments. *Proceedings of the National Academy of Sciences*, page 201714936, January 2018.
- [19] C. Monroe and J. Kim. Scaling the Ion Trap Quantum Processor. *Science*, 339(6124):1164–1169, March 2013.
- [20] M. H. Devoret and R. J. Schoelkopf. Superconducting Circuits for Quantum Information: An Outlook. *Science*, 339(6124):1169–1174, March 2013.
- [21] J. Chiaverini, D. Leibfried, T. Schaetz, M. D. Barrett, R. B. Blakestad, J. Britton, W. M. Itano, J. D. Jost, E. Knill, C. Langer, R. Ozeri, and D. J. Wineland. Realization of quantum error correction. *Nature*, 432(7017):602, December 2004.
- [22] Philipp Schindler, Julio T. Barreiro, Thomas Monz, Volkmar Nebendahl, Daniel Nigg, Michael Chwalla, Markus Hennrich, and Rainer Blatt. Experimental Repetitive Quantum Error Correction. *Science*, 332(6033):1059–1061, May 2011.
- [23] G. Waldherr, Y. Wang, S. Zaiser, M. Jamali, T. Schulte-Herbrüggen, H. Abe, T. Ohshima, J. Isoya, J. F. Du, P. Neumann, and J. Wrachtrup. Quantum error correction in a solid-state hybrid spin register. *Nature*, 506(7487):204, February 2014.
- [24] J. Kelly, R. Barends, A. G. Fowler, A. Megrant, E. Jeffrey, T. C. White, D. Sank, J. Y. Mutus, B. Campbell, Yu Chen, Z. Chen, B. Chiaro, A. Dunsworth, I.-C. Hoi, C. Neill, P. J. J. O’Malley, C. Quintana, P. Roushan, A. Vainsencher, J. Wenner, A. N. Cleland, and John M. Martinis. State preservation by repetitive error detection in a superconducting quantum circuit. *Nature*, 519(7541):66, March 2015.
- [25] M. Veldhorst, C. H. Yang, J. C. C. Hwang, W. Huang, J. P. Dehollain, J. T. Muñonen, S. Simmons, A. Laucht, F. E. Hudson, K. M. Itoh, A. Morello, and A. S. Dzurak. A two-qubit logic gate in silicon. *Nature*, 526(7573):410–414, October 2015.
- [26] Nissim Ofek, Andrei Petrenko, Reinier Heeres, Philip Reinhold, Zaki Leghtas, Brian Vlastakis, Yehan Liu, Luigi Frunzio, S. M. Girvin, L. Jiang, Mazyar Mirrahimi, M. H. Devoret, and R. J. Schoelkopf. Extending the lifetime of a quantum bit with error correction in superconducting circuits. *Nature*, 536(7617):441, July 2016.
- [27] Anton Potocnik, Arno Bargerbos, Florian A. Y. N. Schröder, Saeed A. Khan, Michele C. Collodo, Simone Gasparinetti, Yves Salathé, Celestino Creatore, Christopher Eichler, Hakan E. Türeci, Alex W. Chin, and Andreas Wallraff. Studying Light-Harvesting Models with Superconducting Circuits. *arXiv:1710.07466 [physics, physics:quant-ph]*, October 2017. arXiv: 1710.07466.
- [28] S. Debnath, N. M. Linke, C. Figgatt, K. A. Landsman, K. Wright, and C. Monroe. Demonstration of a small programmable quantum computer with atomic qubits. *Nature*, 536(7614):63–66, August 2016.
- [29] Maika Takita, Andrew W. Cross, A. D. Córcoles, Jerry M. Chow, and Jay M. Gambetta. Experimental Demonstration of Fault-Tolerant State Preparation with Superconducting Qubits. *Physical Review Letters*, 119(18):180501, October 2017.
- [30] T. F. Watson, S. G. J. Philips, E. Kawakami, D. R. Ward, P. Scarlino, M. Veldhorst, D. E. Savage, M. G. Lagally, Mark Friesen, S. N. Coppersmith, M. A. Eriksson, and L. M. K. Vandersypen. A programmable two-qubit quantum processor in silicon. *arXiv:1708.04214 [cond-mat, physics:quant-ph]*, August 2017. arXiv: 1708.04214.
- [31] D. Ristè, C. C. Bultink, K. W. Lehnert, and L. DiCarlo. Feedback Control of a Solid-State Qubit Using High-Fidelity Projective Measurement. *Physical Review Letters*, 109(24):240502, December 2012.
- [32] P. Campagne-Ibarcq, E. Flurin, N. Roch, D. Darson, P. Morfin, M. Mirrahimi, M. H. Devoret, F. Mallet, and B. Huard. Persistent Control of a Superconducting Qubit by Stroboscopic Measurement Feedback. *Physical Review X*, 3(2):021008, May 2013.
- [33] L. Steffen, Y. Salathe, M. Oppiger, P. Kurpiers, M. Baur, C. Lang, C. Eichler, G. Puebla-Hellmann, A. Fedorov, and A. Wallraff. Deterministic quantum teleportation with feed-forward in a solid state system. *Nature*, 500(7462):319–322, August 2013.
- [34] W. Pfaff, B. J. Hensen, H. Bernien, S. B. van Dam, M. S. Blok, T. H. Taminiau, M. J. Tiggelman, R. N. Schouten, M. Markham, D. J. Twitchen, and R. Hanson. Unconditional quantum teleportation between distant solid-state quantum bits. *Science*, 345(6196):532–535, August 2014.
- [35] D. Ristè and L. DiCarlo. Digital feedback in superconducting quantum circuits. In *Superconducting Devices in Quantum Optics*; eds. R.H. Hadfield and G. Johansson; see *arXiv:1508.01385*. Springer, February 2016. arXiv: 1508.01385.
- [36] Navin Khaneja, Timo Reiss, Cindie Kehlet, Thomas Schulte-Herbrüggen, and Steffen J. Glaser. Optimal control of coupled spin dynamics: design of NMR pulse sequences by gradient ascent algorithms. *Journal of Magnetic Resonance*, 172(2):296–305, February 2005.
- [37] S. Machnes, U. Sander, S. J. Glaser, P. de Fouquière, A. Gruslys, S. Schirmer, and T. Schulte-Herbrüggen. Comparing, optimizing, and benchmarking quantum-control algorithms in a unifying programming framework. *Physical Review A*, 84(2):022305, August 2011.
- [38] Peter D. Johnson, Jonathan Romero, Jonathan Olson, Yudong Cao, and Alán Aspuru-Guzik. QVECTOR: an algorithm for device-tailored quantum error correction.

- arXiv:1711.02249 [quant-ph]*, November 2017. arXiv: 1711.02249.
- [39] Alexander Hentschel and Barry C. Sanders. Machine Learning for Precise Quantum Measurement. *Physical Review Letters*, 104(6):063603, February 2010.
- [40] Alexander Hentschel and Barry C. Sanders. Efficient Algorithm for Optimizing Adaptive Quantum Metrology Processes. *Physical Review Letters*, 107(23):233601, November 2011.
- [41] Pantita Palittapongarnpim, Peter Wittek, Ehsan Zahedinejad, Shakib Vedaie, and Barry C. Sanders. Learning in Quantum Control: High-Dimensional Global Optimization for Noisy Quantum Dynamics. *arXiv:1607.03428 [quant-ph, stat]*, July 2016. arXiv: 1607.03428.
- [42] M. Tiersch, E. J. Ganahl, and H. J. Briegel. Adaptive quantum computation in changing environments using projective simulation. *Scientific Reports*, 5:12874, August 2015.
- [43] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. Quantum machine learning. *Nature*, 549(7671):195, September 2017.
- [44] Jonathan Romero, Jonathan P. Olson, and Alan Aspuru-Guzik. Quantum autoencoders for efficient compression of quantum data. *Quantum Science and Technology*, 2(4):045001, 2017.
- [45] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, May 1992.
- [46] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. The MIT Press, Cambridge, Massachusetts, November 2016.
- [47] Peter W. Shor. Scheme for reducing decoherence in quantum computer memory. *Physical Review A*, 52(4):R2493–R2496, October 1995.
- [48] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, Cambridge ; New York, anniversary edition edition, January 2011.
- [49] Barbara M. Terhal. Quantum error correction for quantum memories. *Reviews of Modern Physics*, 87(2):307–346, April 2015.
- [50] Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- [51] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, November 1997.
- [52] Mazyar Mirrahimi, Zaki Leghtas, Victor V. Albert, Steven Touzard, Robert J. Schoelkopf, Liang Jiang, and Michel H. Devoret. Dynamically protected cat-qubits: a new paradigm for universal quantum computation. *New Journal of Physics*, 16(4):045014, 2014.
- [53] Reinier W. Heeres, Philip Reinhold, Nissim Ofek, Luigi Frunzio, Liang Jiang, Michel H. Devoret, and Robert J. Schoelkopf. Implementing a universal gate set on a logical qubit encoded in an oscillator. *Nature Communications*, 8(1):94, July 2017.
- [54] F. Helmer, M. Mariantoni, A. G. Fowler, J. von Delft, E. Solano, and F. Marquardt. Cavity grid for scalable quantum computation with superconducting circuits. *EPL (Europhysics Letters)*, 85(5):50007, 2009.
- [55] R. Li, L. Petit, D. P. Franke, J. P. Dehollain, J. Helsen, M. Steudtner, N. K. Thomas, Z. R. Yoscovits, K. J. Singh, S. Wehner, L. M. K. Vandersypen, J. S. Clarke, and M. Veldhorst. A Crossbar Network for Silicon Quantum Dot Qubits. *arXiv:1711.03807 [cond-mat, physics:quant-ph]*, November 2017. arXiv: 1711.03807.
- [56] Moritz August and José Miguel Hernández-Lobato. Taking gradients through experiments: LSTMs and memory proximal policy optimization for black-box quantum control. *arXiv:1802.04063 [quant-ph]*, February 2018. arXiv: 1802.04063.
- [57] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [58] Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276, February 1998.
- [59] Sham M Kakade. A natural policy gradient. In *Advances in neural information processing systems*, pages 1531–1538, 2002.
- [60] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- [61] Ronald J Williams and Jing Peng. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.
- [62] D Kingma and Jimmy Ba Adam. A method for stochastic optimisation. In *International Conference for Learning Representations*, volume 6, 2015.
- [63] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- [64] Rami Al-Rfou et al. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688, May 2016.

# Supplementary Material

## 1. Physical time evolution

### 1.1. Lindblad equation and completely positive map

We start by discussing how to describe the dynamics of open quantum systems as a brief introduction for people from other fields, and to fix notation.

The state of any Markovian quantum system at time  $t$  is completely characterized by the density matrix  $\hat{\rho}(t)$ . Its time evolution can always be described by a Lindblad equation

$$\frac{d}{dt} \hat{\rho} = -\frac{i}{\hbar} [\hat{H}, \hat{\rho}] + \sum_j \Gamma_j \left( \hat{L}_j \hat{\rho} \hat{L}_j^\dagger - \frac{1}{2} \{ \hat{L}_j^\dagger \hat{L}_j, \hat{\rho} \} \right) \quad (\text{S1})$$

where the Hamiltonian  $\hat{H}$  represents the coherent part of the dynamics and the Lindblad operators  $\hat{L}_j$  (or jump operators) the incoherent part;  $\Gamma_j$  are their corresponding decay rates.  $[\hat{A}, \hat{B}] = \hat{A}\hat{B} - \hat{B}\hat{A}$  is the commutator of two operators, and  $\{\hat{A}, \hat{B}\} = \hat{A}\hat{B} + \hat{B}\hat{A}$  the anti-commutator.  $\hat{A}^\dagger$  denotes the Hermitian conjugate of  $\hat{A}$ .

For compact notation, the terms on the right-hand side can be combined into one superoperator, the Liouvillian  $\mathcal{L}$ :

$$\frac{d}{dt} \hat{\rho} = \mathcal{L} \hat{\rho} \quad (\text{S2})$$

We now introduce the completely positive map  $\Phi(t_f, t_i)$  to formally write down the time evolution of all density matrices in the time interval from  $t_i$  to  $t_f$ :

$$\hat{\rho}(t_f) = \Phi(t_f, t_i)[\hat{\rho}(t_i)] \quad (\text{S3})$$

If  $\mathcal{L}$  does not change over time, the completely positive map is given by  $\Phi(t_f, t_i) = e^{(t_f - t_i)\mathcal{L}}$ ; otherwise, it can be obtained from  $\Phi(t_f, t_i) = \mathcal{T}\exp(\int_{t_i}^{t_f} \mathcal{L}(t) dt)$  where  $\mathcal{T}\exp$  denotes the time-ordered exponential.

Measurements can be nicely integrated into this framework. For each measurement variable and obtained result, there is a projection operator  $\hat{P}$ , and we associate a superoperator  $\mathcal{P}$  defined by its action  $\mathcal{P}\hat{\rho} = \hat{P}\hat{\rho}\hat{P}^\dagger$  on all density matrices  $\hat{\rho}$ . If some measurements are performed at intermediate times  $t_1, \dots, t_m$  with corresponding superoperators  $\mathcal{P}_1, \dots, \mathcal{P}_m$ , and otherwise  $\hat{\rho}$  follows the Lindblad equation characterized by the Liouvillian  $\mathcal{L}$ , the completely positive map has the form  $\Phi(t_f, t_i)[\hat{\rho}] = \phi(t_f, t_i)[\hat{\rho}] / \text{tr}(\phi(t_f, t_i)[\hat{\rho}])$  where  $\phi(t_f, t_i) = e^{(t_f - t_m)\mathcal{L}} \mathcal{P}_m e^{(t_m - t_{m-1})\mathcal{L}} \dots e^{(t_2 - t_1)\mathcal{L}} \mathcal{P}_1 e^{(t_1 - t_i)\mathcal{L}}$  (note that explicit normalization is required because projections are in general not trace-preserving).

For our purposes, it is more convenient to consider  $\phi$  instead of  $\Phi$ , so in the following we will use

$$\hat{\rho}(t_f) = \frac{\phi(t_f, t_i)[\hat{\rho}(t_i)]}{\text{tr}(\phi(t_f, t_i)[\hat{\rho}(t_i)])} \quad (\text{S4})$$

where  $\phi(t_f, t_i)$  is always a linear map (even if it includes measurements).

## 1.2. Simulation of a logical qubit

This formalism is usually applied directly to single density matrices. However, as motivated in the main text, we require an efficient scheme which is capable of processing all possible logical qubit states, i.e. the “full Bloch sphere”, in parallel. This is possible due to two circumstances: the linearity of the completely positive map, and the fact that all initial states are arranged on an affine space.

In Methods, we specify an explicit scheme to construct four quantities  $\hat{\rho}_0(t)$ ,  $\delta\hat{\rho}_x(t)$ ,  $\delta\hat{\rho}_y(t)$  and  $\delta\hat{\rho}_z(t)$  which give access to the density matrix  $\hat{\rho}_{\vec{n}}(t)$  for an arbitrary logical qubit state  $\vec{n} = (x, y, z)$  at any time  $t$ :

$$\begin{aligned}\hat{\rho}_{\vec{n}}(t) &= \frac{\hat{\rho}_0(t) + \sum_j n_j \delta\hat{\rho}_j(t)}{1 + \text{tr}(\sum_j n_j \delta\hat{\rho}_j(t))} = \\ &= \frac{\hat{\rho}_0(t) + x \delta\hat{\rho}_x(t) + y \delta\hat{\rho}_y(t) + z \delta\hat{\rho}_z(t)}{1 + x \text{tr}(\delta\hat{\rho}_x(t)) + y \text{tr}(\delta\hat{\rho}_y(t)) + z \text{tr}(\delta\hat{\rho}_z(t))}\end{aligned}\quad (\text{S5})$$

Without measurements, i.e. under trace-preserving quantum operations like unitary transformations and dissipation, the variables  $\hat{\rho}_0$ ,  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$  follow the same dynamics as the density matrices themselves, i.e. their time evolution is given by the completely positive map  $\Phi(t_f, t_i)$ . However, they have to be treated differently under measurements where renormalization becomes important. Reusing our notion of  $\phi(t_f, t_i)$  from Eq. (S4), we have

$$\begin{aligned}\hat{\rho}_{\vec{n}}(t_f) &= \frac{\phi(t_f, t_i)[\hat{\rho}_{\vec{n}}(t_i)]}{\text{tr}(\phi(t_f, t_i)[\hat{\rho}_{\vec{n}}(t_i)])} = \\ &= \frac{\phi(t_f, t_i)[\hat{\rho}_0(t_i) + \sum_j n_j \delta\hat{\rho}_j(t_i)]}{\text{tr}(\phi(t_f, t_i)[\hat{\rho}_0(t_i) + \sum_j n_j \delta\hat{\rho}_j(t_i)])} = \\ &= \frac{\phi(t_f, t_i)[\hat{\rho}_0(t_i)] + \sum_j n_j \phi(t_f, t_i)[\delta\hat{\rho}_j(t_i)]}{\text{tr}(\phi(t_f, t_i)[\hat{\rho}_0(t_i)] + \sum_j n_j \phi(t_f, t_i)[\delta\hat{\rho}_j(t_i)])}\end{aligned}\quad (\text{S6})$$

We see that it is not *required* to explicitly renormalize  $\hat{\rho}_0$ ,  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$ , and in particular that it is not correct to renormalize them separately (which would be ill-defined because the  $\delta\hat{\rho}_j$  can be traceless). However, a *common* prefactor is allowed as they appear both in the numerator and the denominator, and for convenience we choose the prefactor  $1 / \text{tr}(\phi(t_f, t_i)[\hat{\rho}_0(t_i)])$  such that  $\text{tr}(\hat{\rho}_0(t)) \stackrel{!}{=} 1$  is always satisfied. This leads to the update equations given in Methods.

## 2. Recoverable quantum information

### 2.1. Definition

In the main text, we have defined the recoverable quantum information as

$$\mathcal{R}_Q = \frac{1}{2} \min_{\vec{n}} \|\hat{\rho}_{\vec{n}} - \hat{\rho}_{-\vec{n}}\|_1 \quad (\text{S1})$$

An alternative definition would be

$$\mathcal{R}_Q = \min_{\vec{n}} \left\| \sum_j n_j \delta\hat{\rho}_j \right\|_1 \quad (\text{S2})$$

Both expressions coincide as long as all  $\text{tr}(\delta\hat{\rho}_j) = 0$ , i.e. no information about the logical qubit state has been revealed. For the remaining cases, it is not clear to the authors what is the most natural generalization.

For the following considerations, it will not make a difference which definition is used. Likewise, it does not play a role in the numerical simulations discussed in the main text because, due to our operation sets, only two cases can occur: either no information about the logical qubit has been revealed ( $\frac{1}{2}(\hat{\rho}_{\vec{n}} - \hat{\rho}_{-\vec{n}}) = \sum_j n_j \delta\hat{\rho}_j$ ), or all superpositions have been destroyed (both definitions yield  $\mathcal{R}_Q = 0$ ).

### 2.2. Properties

The purpose of the following discussion is to develop an intuitive understanding for  $\mathcal{R}_Q$ .

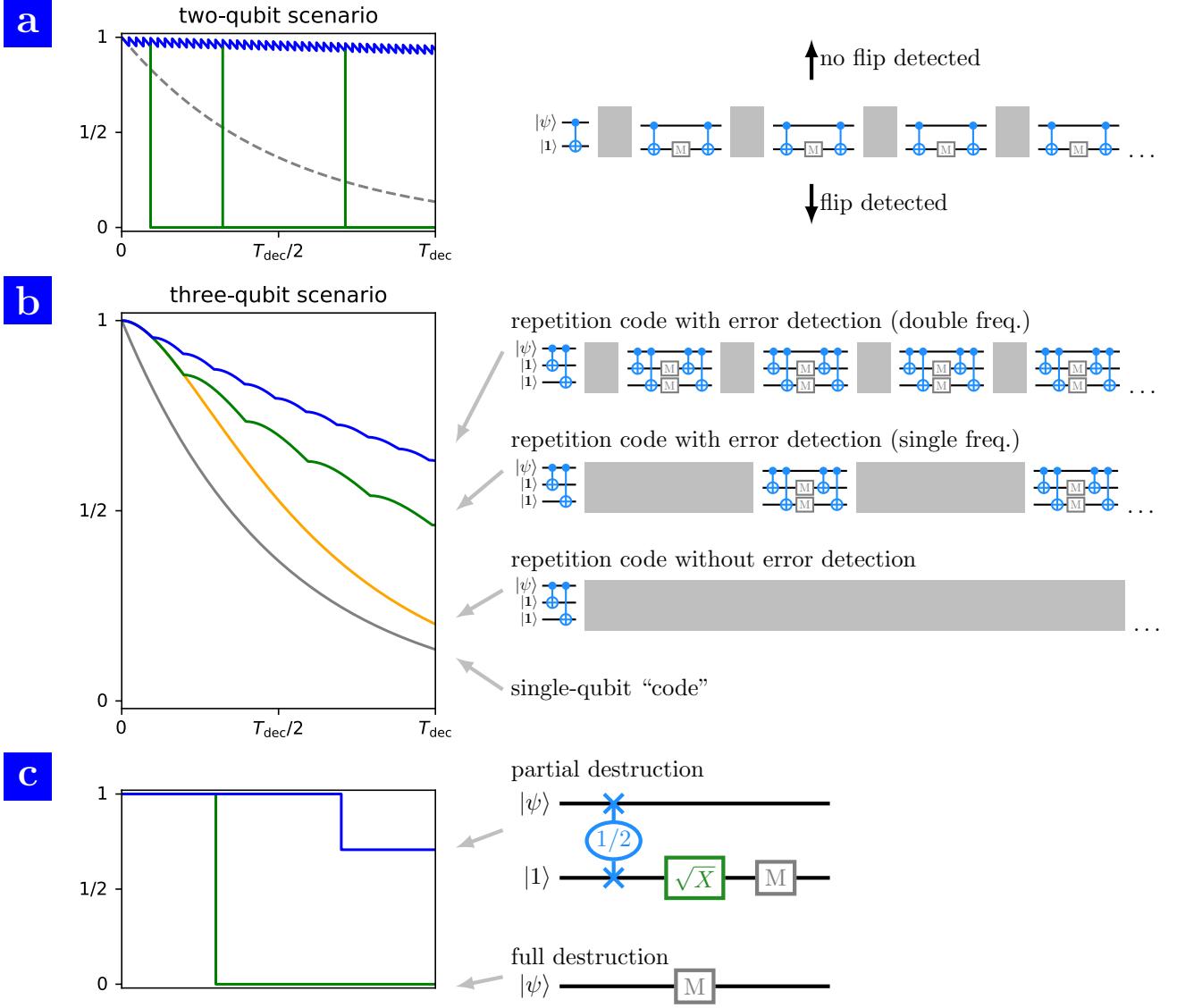


Figure S1. Behavior of  $\mathcal{R}_Q$ . (a) Exemplary curves for  $\mathcal{R}_Q(t)$  in a two-qubit system with bit-flip errors. In the two-qubit encoding  $\alpha|\uparrow\uparrow\rangle + \beta|\downarrow\downarrow\rangle$ , parity measurements can lead to an increase of  $\mathcal{R}_Q$  if no bit flip is detected (blue), but  $\mathcal{R}_Q$  decreases to  $0$  for the alternating result (green lines). The average behavior (gray) is identical to the decay for the single-qubit “encoding”, and for the two-qubit encoding without parity measurements. (b) Decay of  $\mathcal{R}_Q$  in a three-qubit system with bit-flip errors. The curves show the single-qubit “encoding” (gray), the repetition code  $\alpha|\uparrow\uparrow\rangle + \beta|\downarrow\downarrow\rangle$  without error detection (orange) and the repetition code with periodic parity measurements (green and blue). Note that for the two protocols with measurements, we plot the average over all possible measurement results; single trajectories will deviate from this behavior dependent on the particularly found syndroms. In addition, we neglected errors that occur during the syndrom detection sequences, i.e. we implicitly assumed the length of these sequences to be very small compared to the intermediate idle times. This assumption is seriously violated in the bit-flip scenarios on which the networks in the main text are trained (best results for measurements at the period limited by the gate operation time), such that the depicted three-qubit syndrom detection scheme would there perform much worse than the four-qubit schemes found by the neural networks (never enters a fragile state). (c) Drop in  $\mathcal{R}_Q$  as a consequence of total (green) or partial (blue) destruction of superpositions due to a measurement which reveals information about the logical qubit state. The gate sequence leading to the partial drop is  $\sqrt{\text{SWAP}} \rightarrow$  local  $90^\circ$  rotation around  $x$  axis  $\rightarrow$  MSMT( $z$ ).

We start by giving three basic properties of the underlying trace distance  $D_{\vec{n}} = \|\sum_j n_j \delta\hat{\rho}_j(t)\|_1$  for antipodal logical qubit states:

1.  $D_{\vec{n}}$  is invariant under unitary transformations.
2.  $D_{\vec{n}}$  can only decrease under the influence of dissipation (pure states yield the maximum value  $D_{\vec{n}} = 1$ ).
3.  $D_{\vec{n}}$  can be decreased and also increased by measurements, but only under the constraint that the average over all measurement results cannot exceed the prior value:

$$\sum_m p_m D_{\vec{n}}^{(m)} \leq D_{\vec{n}} \quad (\text{S3})$$

Here,  $D_{\vec{n}}$  is the value directly before the measurement, and  $D_{\vec{n}}^{(m)}$  the updated value for result  $m$ ;  $p_m$  denotes the corresponding probability to find this result. (The extreme cases  $\sum_m p_m D_{\vec{n}}^{(m)} = 0$  and  $\sum_m p_m D_{\vec{n}}^{(m)} = D_{\vec{n}}$  can occur and are discussed below).

(1) is a general property of the trace norm. (2) and (3) follow from the contractivity of the trace norm under trace-preserving quantum operations [1, p. 406] (here, the completely positive maps for dissipative dynamics  $\hat{\rho} \mapsto \mathcal{T}\exp(\int_{t_i}^{t_f} \mathcal{L}(t) dt)\hat{\rho}$  and non-selective measurements  $\hat{\rho} \mapsto \sum_m \hat{P}_m \hat{\rho} \hat{P}_m^\dagger$ , respectively). Trace-preserving quantum operations are those which can be written in the form  $\hat{\rho} \mapsto \sum_k \hat{E}_k \hat{\rho} \hat{E}_k^\dagger$  for a complete operator set  $\sum_k \hat{E}_k^\dagger \hat{E}_k = \mathbb{1}$  [1, p. 360].

To illustrate property (3) in an example, we consider the following situation: two qubits are subject to bit-flip errors, i.e.  $\mathcal{D}\hat{\rho} = T_{\text{dec}}^{-1} \sum_q (\hat{\sigma}_z^{(q)} \hat{\rho} \hat{\sigma}_z^{(q)} - \hat{\rho})$ . Initially, the logical qubit is stored in one physical qubit, and the ancilla is prepared in the down state ( $\hat{\sigma}_z |\downarrow\rangle = -|\downarrow\rangle$ ). Then, the protocol depicted in Fig. S1a is applied: first, a CNOT gate entangles the data qubit with the ancilla. During the following idle time,  $\mathcal{R}_Q$  decays like  $\mathcal{R}_Q \sim e^{-2t/T_{\text{dec}}}$ . Afterwards, the qubits are disentangled again by a second CNOT gate (which leaves  $\mathcal{R}_Q$  invariant). If in this situation a measurement on the ancilla reveals that it is still in the down state,  $\mathcal{R}_Q$  grows to  $\mathcal{R}_Q^{(\downarrow)} = 2e^{-2t/T_{\text{dec}}} / (1 + e^{-4t/T_{\text{dec}}}) > e^{-2t/T_{\text{dec}}}$ . However, the probability to find the ancilla in the down state is only  $P_\downarrow = \frac{1}{2}(1 + e^{-4t/T_{\text{dec}}})$ , and if a flip is detected ( $P_\uparrow = 1 - P_\downarrow$ ),  $\mathcal{R}_Q$  drops to 0 as it cannot be determined which qubit has flipped. Hence, *on average* nothing is won:  $P_\downarrow \mathcal{R}_Q^{(\downarrow)} + P_\uparrow \mathcal{R}_Q^{(\uparrow)} = \mathcal{R}_Q$ . (This basically shows that the described protocol allows quite well to detect errors, but it is not possible to resolve them, and thus it cannot be used to slowdown the decay.) Exemplary  $\mathcal{R}_Q$  curves for the repeated application of this protocol are shown in Fig. S1a.

We proceed by considering the textbook three-qubit scenario with bit-flip error channels [1, p. 427ff] and inspecting it from the perspective of  $\mathcal{R}_Q$  (see also Fig. S1b). Compared to the decay  $\mathcal{R}_Q(t) = e^{-2t/T_{\text{dec}}}$  for the “trivial” single-qubit encoding, e.g.  $(\alpha|\uparrow\rangle + \beta|\downarrow\rangle) \otimes |\downarrow\downarrow\rangle$ , the encoding  $\alpha|\uparrow\uparrow\uparrow\rangle + \beta|\downarrow\downarrow\downarrow\rangle$  (bit-flip code, repetition code) improves the situation in the beginning:  $\mathcal{R}_Q(t) = \frac{1}{2}(3 - e^{-4t/T_{\text{dec}}})e^{-2t/T_{\text{dec}}}$ . However, the asymptotical behavior for  $t \rightarrow \infty$  is still the same. This behavior can be explained as follows: as the populations for single bit-flips accumulate over time, the probability grows that a second bit-flip (on another qubit) occurs which cannot be resolved anymore. The only way to actually protect the logical qubit is a suitable error detection scheme, i.e. to keep track of possible bit-flips via repeated parity measurements. Such a measurement delays the decay again for some time: after a syndrom detection at  $t_0$ , we have  $\mathcal{R}_Q(t > t_0) = \frac{1}{4}(3 - e^{-4t_0/T_{\text{dec}}})(3 - e^{-4(t-t_0)/T_{\text{dec}}})e^{-2t/T_{\text{dec}}}$  (assuming a quasi-instantaneous syndrom detection during which no further errors can occur). Thus, even though measurements cannot (on average) improve  $\mathcal{R}_Q$  immediately (cmp. property 3), they are useful (and necessary) to prevent future losses in  $\mathcal{R}_Q$ .

Last, but not least,  $\mathcal{R}_Q$  correctly describes the loss in recoverability of the quantum state due to measurements which reveal information about the logical qubit state (see also Fig. S1c). If all superpositions are destroyed,  $\mathcal{R}_Q$  drops to 0 (independent of the measurement result); this happens, for example, in the “trivial encoding” (with the logical qubit stored directly in one physical qubit) if a measurement is performed on this data qubit. In comparison, partial destruction of superpositions is reflected by an accordingly smaller decrease in  $\mathcal{R}_Q$ . By contrast, measurements which reveal no information about the logical qubit state do not change  $\mathcal{R}_Q$  *on average* ( $\sum_m p_m D_{\vec{n}}^{(m)} = D_{\vec{n}}$ ).

### 2.3. Computation

The numerical evaluation of  $\mathcal{R}_Q$  is technically challenging. The computation of the trace norm  $\|\sum_j n_j \delta\hat{\rho}_j(t)\|_1$  for a given point  $\vec{n}$  on the Bloch sphere is already expensive as it involves the eigendecomposition of its argument. Moreover, the dependence of  $\|\sum_j n_j \delta\hat{\rho}_j(t)\|_1$  on the Bloch vector  $\vec{n}$  can be complex, which makes it non-trivial to locate the minimum that determines  $\mathcal{R}_Q$ .

In our experience, it is not valid to approximate the trace-norm  $\|\cdot\|_1$  in  $\mathcal{R}_Q$  using the Hilbert-Schmidt norm  $\|\cdot\|_2$ , which is numerically far easier to calculate (and minimize). Likewise, other simpler quantities (like the purity of the multi-qubit state) are not sufficient to replace  $\mathcal{R}_Q$ .

In the Methods, we have claimed that for CHZ circuits (all unitary operations are combinations of CNOT, Hadamard and phase gates; measurement of variables in the Pauli group; state preparation in the computational basis [1, p. 464]), the anti-commutator relation  $\{\delta\hat{\rho}_j, \delta\hat{\rho}_k\} = 0$  is satisfied for all distinct  $j, k \in \{x, y, z\}$ ; our bit-flip scenarios fall into this category of CHZ circuits. Furthermore, we have argued that  $\min_{x,y,z} \|x\delta\hat{\rho}_x + y\delta\hat{\rho}_y + z\delta\hat{\rho}_z\|_1 = \min_{j \in \{x,y,z\}} \|\delta\hat{\rho}_j\|_1$  (for  $x^2 + y^2 + z^2 = 1$ ) if this anti-commutator relation is satisfied. In the following, we briefly sketch the proof for this statement. We note that “ $\leq$ ” is a priori clear, so just “ $\geq$ ” has to be shown. The assumption  $j \neq k \Rightarrow \{\delta\hat{\rho}_j, \delta\hat{\rho}_k\} = 0$  directly implies that  $[\delta\hat{\rho}_j^2, \delta\hat{\rho}_k^2] = 0$ , i.e. there is a common eigenbasis  $\{|\phi_n\rangle\}$  for all three  $\delta\hat{\rho}_j^2$ . Further, we can conclude that

$$\|x\delta\hat{\rho}_x + y\delta\hat{\rho}_y + z\delta\hat{\rho}_z\|_1 = \sum_n \sqrt{x^2\langle\phi_n|\delta\hat{\rho}_x^2|\phi_n\rangle + y^2\langle\phi_n|\delta\hat{\rho}_y^2|\phi_n\rangle + z^2\langle\phi_n|\delta\hat{\rho}_z^2|\phi_n\rangle} \quad (\text{S4})$$

and for  $x^2 + y^2 + z^2 = 1$ , Jensen’s inequality (note that square root is a concave function) yields the lower bound

$$\begin{aligned} \|x\delta\hat{\rho}_x + y\delta\hat{\rho}_y + z\delta\hat{\rho}_z\|_1 &\geq \sum_n x^2 \sqrt{\langle\phi_n|\delta\hat{\rho}_x^2|\phi_n\rangle} + y^2 \sqrt{\langle\phi_n|\delta\hat{\rho}_y^2|\phi_n\rangle} + z^2 \sqrt{\langle\phi_n|\delta\hat{\rho}_z^2|\phi_n\rangle} = \\ &= x^2\|\delta\hat{\rho}_x\|_1 + y^2\|\delta\hat{\rho}_y\|_1 + z^2\|\delta\hat{\rho}_z\|_1 \end{aligned} \quad (\text{S5})$$

Minimization over the Bloch sphere ( $x^2 + y^2 + z^2 = 1$ ) leads to the desired result.

### 3. State-Aware Network

#### 3.1. Reward scheme

*Protection reward* The purpose of the protection reward is that the agent learns how to protect the quantum information against the error mechanisms of the quantum system, i.e. to preserve a state that could in principle be recovered without the need to actually perform the recovery sequence (also cmp. Methods). For instance, in the bit-flip examples, the agent learns to go into a proper encoding and to perform afterwards repeated parity measurements. Most of the networks in sec. “Results” of the maintext are trained solely with this type of reward; only for fig. 7 in the main text, an additional recovery reward (see below) has been considered.

For convenience, we will set the simulation time step to  $\Delta t = 1$  in this discussion.

As discussed in appendix 2 and the main text,  $\mathcal{R}_Q$  is a powerful measure for the capability to recover a target state (without the need to actually perform these steps), and so the objective is to maximize the value of  $\mathcal{R}_Q$  at the final time  $T$ . For learning efficiency, and (very important for us) to credit already partial success, we aim for an immediate reward. This reward should essentially be of the form  $\tilde{r}_t = \mathcal{R}_Q(t+1) - \mathcal{R}_Q(t)$ : as the learning algorithm optimizes the reward sum,  $\sum_t \tilde{r}_t = \sum_t \mathcal{R}_Q(t+1) - \mathcal{R}_Q(t) = \mathcal{R}_Q(T) - \mathcal{R}_Q(0) = \mathcal{R}_Q(T) - 1$ , the agent is guided towards the desired behavior, to maximize  $\mathcal{R}_Q(T)$ .

In practice, we have found it useful to implement the reward scheme rather in the following form, which will be motivated below:

$$r_t = \begin{cases} 1 + \frac{\bar{\mathcal{R}}_Q(t+1) - \mathcal{R}_Q(t)}{2\Delta t/T_{\text{triv}}} + 0 & \text{if } \bar{\mathcal{R}}_Q(t+1) > 0 \\ 0 & \text{if } \bar{\mathcal{R}}_Q(t) \neq 0 \wedge \bar{\mathcal{R}}_Q(t+1) = 0 \\ 0 & \text{if } \bar{\mathcal{R}}_Q(t) = 0 \end{cases} \quad (\text{S1})$$

$$=: r_t^{(1)} \quad =: r_t^{(2)}$$

where  $\Delta t$  denotes the time per gate operation,  $T_{\text{triv}}$  is the decay time of the quantum state in the trivial encoding, and  $P$  is the punishment for a measurement which reveals the logical qubit state.  $\bar{\mathcal{R}}_Q$  will be defined in the next paragraph. Note that  $r_t$  is still designed such that optimizing  $\mathcal{R}_Q(T)$  leads to the highest possible reward sum  $\sum_t r_t$ .

The updated value of  $\mathcal{R}_Q$  after a measurement can depend on its outcome. We define  $\bar{\mathcal{R}}_Q(t+1)$  as the expectation value for  $\mathcal{R}_Q(t+1)$ . For measurements, this is the average for all possible measurement results, weighted according to their probability. For unitary gates, there is no ambiguity and thus  $\bar{\mathcal{R}}_Q(t+1) = \mathcal{R}_Q(t+1)$ . We will motivate below why it makes sense to consider this quantity.

When computing the (discounted) return  $R_t$  from the reward, we choose to distribute only  $r_t^{(1)}$  backwards in time, while  $r_t^{(2)}$  is assigned directly to the respective action:

$$R_t = (1 - \gamma) \left( \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}^{(1)} \right) + r_t^{(2)} \quad (\text{S2})$$

$\gamma$  is the discount rate; we choose a value of  $\gamma = 0.95$ , corresponding to a half-life in the range between 10 and 20 time steps.

The expressions in Eqs. (S1) and (S2) are chosen for the following reasons:

- For measurements, there is the special situation that, dependent on the quantum state, they might reveal information about the logical qubit state. This is always accompanied by the destruction of superpositions, and in our operation sets, it even leads to their complete destruction. Because no other action could lead to a worse state and the exceptionally strong variation in  $\mathcal{R}_Q$  would cause learning instabilities, it makes sense to treat those cases separately.

These “destructive” measurements are characterized by the fact that  $\mathcal{R}_Q$  drops to 0 independent of the measurement result. If this situation occurs, we deviate from a reward that is directly related to  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t)$ , and instead set the immediate reward  $r_t$  to a fixed negative value  $-P$ . Furthermore, because this situation could easily be avoided by choosing this particular action different, we do not distribute this reward over the previous time steps when we compute the return.

As an alternative, those measurements could simply be excluded from being performed by an external instance. We decided against doing so because in the real-world application we aim for (control of a quantum experiment), an agent needs to be able to detect these cases on its own.

- Even after excluding those measurements which reveal the logical qubit state,  $\mathcal{R}_Q$  can still jump in special situations; for example, in the two-qubit scenario discussed in appendix 2.2, there are measurements that make  $\mathcal{R}_Q$  increase in case of the anticipated result, at the expense of getting  $\mathcal{R}_Q = 0$  for the unlikely one. These strong fluctuations in  $\mathcal{R}_Q$  would lead to large variations in the reward and thereby can easily cause learning instabilities. This can be improved by using the variable  $\bar{\mathcal{R}}_Q$  as defined above: we can easily replace  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t) \rightarrow \bar{\mathcal{R}}_Q(t+1) - \mathcal{R}_Q(t)$  because in the end, the reward is used to compute an estimator (for the learning gradient, see appendix 3.3), and on average both coincide. Since  $\bar{\mathcal{R}}_Q$  is much more stable than  $\mathcal{R}_Q$ , we can easily get rid of the instabilities; as an intended side-effect, this also reduces the “normal” noise level on the learning gradient (due to the deviations between the estimator and the true value) because now each simulation averages directly over all possible measurement results, something that otherwise would have to be done over many epochs.
- The first case in Eq. (S1) is a scaled version of  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t)$  due to the following consideration: For the trivial encoding, we have (at the start of the simulation)  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t) \approx -2\Delta t/T_{\text{triv}}$  since in this case  $\mathcal{R}_Q(t) \approx e^{-2t/T_{\text{triv}}} \approx 1 - 2t/T_{\text{triv}}$  for  $t \ll T_{\text{triv}}$ , whereas for perfect conservation of the quantum information, i. e. constant  $\mathcal{R}_Q(t)$ , we would get  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t) = 0$ . Note that typically, even the ideal action sequence does not reach constant  $\mathcal{R}_Q(t)$ , but it provides a good approximation for comparison with other strategies. From comparing the values 0 and  $-2\Delta t/T_{\text{triv}}$ , we can see that the typical value range for  $\tilde{r}_t = \mathcal{R}_Q(t+1) - \mathcal{R}_Q(t) = 0$  depends on parameters of the physical model. In practice, this makes it difficult to change these properties, especially the type of error channels and the corresponding decoherence rates. For a proper normalization, we can simply divide  $\mathcal{R}_Q(t+1) - \mathcal{R}_Q(t)$  by  $2\Delta t/T_{\text{triv}}$ ; in addition, we can always add the constant value which we choose to be 1, such that the trivial encoding earns  $r_t \approx 0$  and the ideal strategy  $r_t \approx 1$ .

*Recovery reward* In practice, the goal is to eventually recover the logical qubit state from the physical qubits. This means that the quantum system should finally be brought into a state where the logical qubit state can be read off easily from one target qubit, with a pre-defined interpretation for the polarization. To train the network towards this behavior, we will introduce a recovery reward  $r_t^{(\text{recov})}$  (see below) that is given in addition to the protection reward (which only addresses preserving the recoverable quantum information  $\mathcal{R}_Q$ ). Furthermore, we extend the input of the neural network to also contain the time. It is fully sufficient to add a “countdown” for the last time steps to signal the network when the decoding should start (when also the decoding reward can be earned). We input this countdown in a one-hot encoding. Note that it does not matter how many time steps exactly remain when giving the decoding signal (and the recovery reward), if there are sufficiently many left to fit the decoding sequence, which of course is initially unknown. For our example shown in fig. 7 of the main text the countdown is started when  $t > T - 10$ , and the decoding reward is only given afterwards.

A successful recovery requires to correct for the errors (if some occurred) and to decode, such that in the end the logical qubit is stored solely on a specified target qubit. All other qubits should be disentangled and in particular the state of the target qubit should not be flipped compared to the initial state. To achieve this, we extend the reward scheme by two contributions: First, an additional decoding reward can be earned during the last few time steps. Second, a correction reward can be earned at the last time step (but only if the decoding has been successfully performed). We combine them into the recovery reward

$$r_t^{(\text{recov})} = \underbrace{\beta_{\text{dec}} (D_{t+1} - D_t)}_{\text{decoding reward}} + \underbrace{\beta_{\text{corr}} C_T \delta_{t,T-1}}_{\text{correction reward}} \quad (\text{S3})$$

where  $\beta_{\text{dec}}$  and  $\beta_{\text{corr}}$  are the scalings for the decoding and correction reward respectively. In the following, we will discuss first the explicit form of  $D_t$  (for decoding) and later that of  $C_t$  (for correction).

We design the decoding reward to ease learning for the network. Instead of providing a reward only if all qubits of our example are in their respective desired target state of either containing or not containing information about the logical qubit state, it helps to immediately reward also partial success, i. e. whenever one more qubit is brought into the desired state. Thereby, this disentanglement of qubits can be learned stepwise until all qubits are decoded correctly.

In order to quantify such a reward, we aim for a criterion to determine whether a qubit (ignoring all other qubits) contains at least a small fraction of quantum information about the logical qubit state. We obtain a mathematically sufficient criterion by investigating the partial trace over all other qubits: in order to conclude that qubit  $q$  definitely contains information about the logical qubit state, we evaluate whether

$$\text{tr}_{\bar{q}}(\hat{\rho}_{x,y,z}(t)) \quad (\text{S4})$$

depends on the logical qubit state  $\vec{n} = (x, y, z)$ . Expressed in terms of the variables  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$  (cmp. Eq. (S5)), stating that the quantity  $\text{tr}_{\bar{q}}(\hat{\rho}_{x,y,z}(t))$  does depend on  $\vec{n}$  and therefore  $q$  carries information about the logical qubit state is equivalent to the condition that at least one of the corresponding partial traces is non-vanishing:

$$(\text{tr}_{\bar{q}}(\delta\hat{\rho}_x(t)) \neq 0) \vee (\text{tr}_{\bar{q}}(\delta\hat{\rho}_y(t)) \neq 0) \vee (\text{tr}_{\bar{q}}(\delta\hat{\rho}_z(t)) \neq 0) \quad (\text{S5})$$

We have successfully applied this criterion to train the network analyzed in fig. 7 of the main text. To see why this is strictly speaking not a necessary condition (and thus not generally applicable), consider the Laflamme-Miquel-Paz-Zurek [2] encoding. There, this criterion would predict that there is no quantum information in any of the qubits, even though the expected result is that it is distributed over all of them. Nonetheless, for our example of the repetition code it is very useful to help the network to learn the decoding with a close to perfect success rate. One can now specify a desired target state, where e.g. the first qubit is supposed to contain the quantum information, but none of the remaining qubits. By evaluating condition Eq. (S5) for each qubit, we identify how many qubits are in the desired target state of containing or not-containing quantum information and calculate

$$D_t = \begin{cases} 0 & \text{if } t < T_{\text{signal}} \\ \frac{\#\text{correct qubits} - \#\text{wrong qubits}}{2} & \text{if } t \geq T_{\text{signal}} \end{cases} \quad (\text{S6})$$

where  $T_{\text{signal}}$  denotes the time where the neural net obtains the signal to decode, i. e., where the countdown starts. As the decoding reward is essentially  $D_{t+1} - D_t$  (see Eq. (S3)), a positive reward is given if the number of qubits in the desired state increases from one time step to the next, and the network is punished with a negative reward if the number of qubits in the desired state decreases. The corresponding coefficient  $\beta_{\text{dec}}$  has to be sufficiently large such that the decoding reward can compete with the protection reward (we choose  $\beta_{\text{dec}} = 20$ ). This is because during and after the decoding,  $\mathcal{R}_Q$  decays faster, lowering the protection reward. Note that the decoding reward is set to zero if the recoverable quantum information is already too small. In our particular example, we have chosen the threshold to be  $\mathcal{R}_Q < 0.1$ .

Having decoded the multi-qubit state leaves the logical qubit solely on a specified physical qubit. However, this is not necessarily the initial logical qubit state. Preserving the recoverable quantum information  $\mathcal{R}_Q$  only implies that the initial logical qubit state could in principle be recovered from the actual qubit state. To trigger the network to perform corrections, such that the final state is not flipped with respect to the initial logical qubit state, we calculate whether the Bloch vector of the specified physical qubit after the final time step is rotated with respect to the Bloch vector of the initial logical qubit state. To quantify this, we consider the overlap  $\langle \phi_{\vec{n}} | \text{tr}_{\text{o.q.}}(\hat{\rho}_{\vec{n}}) | \phi_{\vec{n}} \rangle$  for a given logical qubit state  $\vec{n}$ ;  $\text{tr}_{\text{o.q.}}$  denotes the partial trace over all qubits except for the target qubit. As a successful error correction scheme needs to work for all possible logical qubit states, we focus on the worst-case, i. e. we consider

$$\mathcal{O}_Q = \min_{\vec{n}} \langle \phi_{\vec{n}} | \text{tr}_{\text{o.q.}}(\hat{\rho}_{\vec{n}}) | \phi_{\vec{n}} \rangle \quad (\text{S7})$$

(cmp. Eq. (S6)). According to *Eq. (S5)*, we can write  $\hat{\rho}_{\vec{n}}$  in terms of  $\hat{\rho}_0$ ,  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$ . If  $\text{tr}_{\text{o.q.}}(\hat{\rho}_0) = \frac{1}{2}\mathbb{1}$  and no information about the logical qubit state has been revealed ( $\text{tr}_{\text{o.q.}}(\delta\hat{\rho}_j)$  for all  $j \in \{x, y, z\}$ ),  $\mathcal{O}_Q$  simplifies to

$$\mathcal{O}_Q = \frac{1 + \text{mineig } \hat{A}}{2} \quad (\text{S8})$$

with  $A_{jk} = \frac{1}{4}(\langle \phi_{\vec{e}_j} | \text{tr}_{\text{o.q.}}(\delta\hat{\rho}_k) | \phi_{\vec{e}_j} \rangle + \langle \phi_{\vec{e}_k} | \text{tr}_{\text{o.q.}}(\delta\hat{\rho}_j) | \phi_{\vec{e}_k} \rangle)$ ;  $\text{mineig } \hat{A}$  denotes the smallest eigenvalue. In our bit-flip scenarios, only the two cases  $\text{mineig } \hat{A} = \pm \mathcal{R}_Q$  can occur after decoding into the target qubit. We set  $C_T = 1$  if (at the end of the simulation) the decoding has been performed completely and  $\text{mineig } \hat{A} > 0$ , and otherwise  $C_T = 0$ . Like for the decoding reward, the coefficient  $\beta_{\text{corr}}$  has to be sufficiently large such that the correction reward can compete with the protection reward (we choose  $\beta_{\text{corr}} = 10$ ).

In the learning gradient, we do not consider directly the reward  $r_t^{(\text{recov})}$ , but rather the return which we choose as

$$R_t^{(\text{recov})} = (1 - \gamma) \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}^{(\text{recov})}. \quad (\text{S9})$$

with discount rate  $\gamma$ .

### 3.2. Network Input

As discussed in Methods, the input to the state-aware network consists of (i) a representation of four evolved density matrix contributions (that represent the evolution of all logical qubit states), (ii) one neuron per measurement in the action set which predicts whether executing that particular measurement would reveal information about the logical qubit state, (iii) the previous action, and (iv) a counter indicating the physical time during the last few time steps; (iv) is given only if we want the agent to perform explicit recovery. Here, we provide additional information for (i) and (ii):

- As the state-aware network is supposed to have perfect information, its input should make it possible to reconstruct the quantum state for any logical qubit state.

This information is contained in  $\mathcal{A} = \{\hat{\rho}_0, \delta\hat{\rho}_x, \delta\hat{\rho}_y, \delta\hat{\rho}_z\}$ , or any (non-trivial) combination of these quantities. In particular, this includes  $\mathcal{B} = \{\hat{\rho}_0, \hat{\rho}_1 := \hat{\rho}_0 + \delta\hat{\rho}_x, \hat{\rho}_2 := \hat{\rho}_0 + \delta\hat{\rho}_y, \hat{\rho}_3 := \hat{\rho}_0 + \delta\hat{\rho}_z\}$ ;  $\mathcal{B}$  consists of the density matrices for the logical qubit states  $\mathbb{1}/2$  (fully depolarized),  $\hat{\sigma}_x$ ,  $\hat{\sigma}_y$  and  $\hat{\sigma}_z$  (up to normalization which only matters if information about the logical qubit state has been revealed). We prefer  $\mathcal{B}$  over  $\mathcal{A}$  because it is in a more “compact” format:  $\hat{\rho}_1$ ,  $\hat{\rho}_2$  and  $\hat{\rho}_3$  are usually dominated by one eigenstate (as long as the quantum information is well preserved), whereas  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$  already start from two components with equally large contribution, so the double number of PCA components would be required for the same amount of information. In practice, we have found that this indeed makes a difference in learning.

The straightforward approach would be to feed these density matrices directly into the neural network (one input neuron per matrix component). However, we have found it useful to pre-process the density matrices via principal component analysis (PCA). The key motivation is to reduce the input size of the neural network: compared to the  $4^{\#\text{qubits}}$  entries of a density matrix, we downsize the input to  $P(\#\text{qubits}) \cdot 2^{\#\text{qubits}}$  components where  $P(\#\text{qubits})$  is the number of PCA components which we would typically choose polynomially in  $\#\text{qubits}$  (note that the eigenstates are of dimension  $2^{\#\text{qubits}}$ , and the density matrices of dimension  $2^{\#\text{qubits}} \times 2^{\#\text{qubits}}$ ). So, the input size still grows exponentially with  $\#\text{qubits}$  even for the PCA-ed input, but compared to the non-PCA case we still win an exponential factor. A nice side-effect is that the PCA automatically decomposes density matrices into main channels and error contributions (and sorts these components by influence); we suppose that this helps the network to better recognize the input states.

- As already explained in appendix 3.1, a particular pitfall are measurements which reveal information about the logical qubit and thereby “destroy” the quantum state. In the following, we describe an extension to the input which helps the network to detect these cases.

From Eq. (S5), it can be seen that the knowledge about the logical qubit state is determined by the vector

$$\vec{b} := \begin{pmatrix} \text{tr}(\delta\hat{\rho}_x) \\ \text{tr}(\delta\hat{\rho}_y) \\ \text{tr}(\delta\hat{\rho}_z) \end{pmatrix} \quad (\text{S10})$$

The ideal case where no information has leaked out, i.e. all logical states are still equally likely, is equivalent to  $\vec{b} = \vec{0}$ . We now consider the effect of measurements on  $\vec{b}$ . For a measurement with two possible results (let  $\hat{P}_+$  and  $\hat{P}_-$  be the corresponding projection operators), the “measurement bias” vector

$$\Delta\vec{b} := \begin{pmatrix} \text{tr}((\hat{P}_+ - \hat{P}_-) \delta\hat{\rho}_x) \\ \text{tr}((\hat{P}_+ - \hat{P}_-) \delta\hat{\rho}_y) \\ \text{tr}((\hat{P}_+ - \hat{P}_-) \delta\hat{\rho}_z) \end{pmatrix} \quad (\text{S11})$$

describes the change of  $\vec{b}$ :

$$\vec{b} \xrightarrow{\pm} \frac{\vec{b} \pm \Delta\vec{b}}{2 \text{tr}(\hat{P}_\pm \hat{\rho}_0)} \quad (\text{S12})$$

Therefore,  $\Delta\vec{b}$  indicates how much additional information about the logical qubit state is gathered by performing the corresponding measurement.

For our set of operations, each measurement is either unbiased or leads to a complete “collapse” of the quantum state; from Eq. (S12), we see that the condition  $\Delta\vec{b} \stackrel{?}{=} \vec{0}$  should distinguish these cases. So, we compute  $\Delta\vec{b} \stackrel{?}{=} \vec{0}$  for each measurement in the action set and provide it as additional input to the network. Because this information can be extracted from  $\delta\hat{\rho}_x$ ,  $\delta\hat{\rho}_y$  and  $\delta\hat{\rho}_z$ , this is in the strict sense no additional input, but rather a “rewording” of a special property in a convenient format.

We have observed that these neurons indeed accelerate the learning process, especially in the early training phase.

### 3.3. Reinforcement Learning Algorithm

Reinforcement learning[3] is a general framework to autonomously explore strategies for “optimum control” problems. In the terminology of reinforcement learning, the control problem is represented by an “environement” and the controller by an “agent”. This agent can successively choose between different actions; for this purpose, it typically has (at least partial) knowledge about state of the environment.

In the policy gradient[4] approach, the agent directly computes a policy function  $\pi_\theta(a_t|s_t)$  which gives the probability to choose action  $a_t$  in state  $s_t$ ;  $\theta$  is a (multi-dimensional) parameter representing the internal state of the agent, for us typically the weights and biases of the neural network. Training means to search for weights  $\theta$  which yield some desired behavior.

The simplest way to train such a policy is given by the “vanilla” policy gradient

$$g_{\text{van}} = \mathbb{E} \left[ \sum_t R_t \frac{\partial}{\partial \theta} \ln \pi_\theta(a_t|s_t) \right] \quad (\text{S13})$$

where  $\mathbb{E}$  denotes the expectation value (according to the policy  $\pi_\theta(a_t|s_t)$ ) and  $R_t$  the return collected for the particular action sequence (see appendix 3.1 for our choice of the return). In practice, an estimate for  $g_{\text{van}}$  is computed from a finite number of simulations, and then the parameters  $\theta$  are updated – again in the simplest case – according to the update rule

$$\theta \rightarrow \theta + \eta g_{\text{van}} \quad (\text{S14})$$

(“steepest ascent”) with the learning rate  $\eta$ ; this procedure defines one epoch, and is repeated multiple times.

There are various modifications to the learning scheme as discussed so far which often lead to significant improvements in terms of learning speed and stability. Specifically, we make use of the following techniques:

- Instead of Eq. (S14), we actually use the adaptive moment estimation (Adam) update rule[5]

$$\theta \rightarrow \theta + \eta \cdot B_n \cdot \frac{m}{\sqrt{v}} \quad (\text{S15a})$$

$$m \rightarrow \beta_1 m + (1 - \beta_1) g \quad (\text{S15b})$$

$$v \rightarrow \beta_2 v + (1 - \beta_2) g^{\odot 2} \quad (\text{S15c})$$

with the Adam hyperparameters  $\eta$ ,  $\beta_1$  and  $\beta_2$ ;  $\odot^2$  denotes the element-wise square, and with  $\sqrt{v}$  we mean the element-wise square root of  $v$ . The coefficient  $B_n$  can depend on the epoch index  $n$  and is proposed by the Adam

inventors to counteract the zero bias for  $m$  and  $v$  in the early training phase ( $B_n = \sqrt{1 - \beta_2^n}/(1 - \beta_1^n)$ ); in the results shown here and in the main text, we did not employ bias correction which is realized by  $B_n = 1$ . Both variants yield comparable results for us.

- A particular difficulty which occurs in our challenge is that we have a quasi-continuous spectrum for the reward, and it is necessary to resolve small differences in there. Improving this aspect is one of the advantages of the natural policy gradient[6], defined by

$$Fg_{\text{nat}} = g_{\text{van}} \quad (\text{S16})$$

with the Fisher information matrix

$$F = \mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \ln \pi_\theta(a|s) \right) \left( \frac{\partial}{\partial \theta} \ln \pi_\theta(a|s) \right)^\top \right] \quad (\text{S17})$$

for actions  $a$  and states  $s$  being distributed according to the policy  $\pi_\theta$ ;  $^\top$  denotes vector transpose. For more information, see [6–10] (the appendix of the latter also gives a recipe for an efficient implementation based on L and R operations).

- Policy gradient schemes can often be enhanced by subtracting a baseline from the return, i. e.

$$g_{\text{van}} \rightarrow \mathbb{E} \left[ \sum_t (R_t - b) \frac{\partial}{\partial \theta} \ln \pi_\theta(a_t|s_t) \right] \quad (\text{S18})$$

where  $b$  is in the simplest case an average of past values for the return. The motivation is that for suitable  $b$ , the variance of  $R_t - b$  is smaller than the variance of  $R_t$  alone, and so the estimator for this modified learning gradient should have smaller fluctuations.

There is a tailored way to choose the baseline  $b$  for the natural policy gradient [9], but we do not use this one because another aspect is more important for us. Due to the structure of the protection reward (see appendix 3.1), the corresponding return is explicitly time-dependent: for successful strategies where  $r_t^{(1)} \approx 1$  and  $r_t^{(2)} \approx 0$ , the return behaves as  $R_t \approx (1 - \gamma) \sum_{k=0}^{T-t-1} \gamma^k = 1 - \gamma^{T-t}$ , i. e.  $R_t$  drops from 1 to 0 near the end of a trajectory. Also the distribution of the recovery reward in time is highly uneven (only assigned in the final time steps). To compensate for this time dependence of the return, we use a time-dependent baseline:

$$g = \mathbb{E} \left[ \sum_t (R_t - b_t) \frac{\partial}{\partial \theta} \ln \pi_\theta(a_t|s_t) \right] \quad (\text{S19})$$

In practice, we choose the baseline  $b_t$  as exponentially decaying average of the return  $R_t$ : in epoch  $N$ , the baseline takes the value

$$b_t = (1 - \kappa) \sum_{n=0}^{N-1} \kappa^n \bar{R}_t^{(N-1-n)} \quad (\text{S20})$$

where  $\kappa = 0.9$  is the discount rate and  $\bar{R}_t^{(n)}$  is the mean of the return at time step  $t$  found in epoch  $n$ .

We have observed that this explicitly time-dependent baseline considerably improves learning stability and leads to much faster learning.

- Without any countermeasures, policies often strongly tend to prefer few actions only. In order to encourage exploration, it has been found useful to introduce entropy regularization [11]: the learning gradient  $g$  used in the update equation (cmp. Eqs. (S14) and (S15)) is substituted by  $g + \lambda \partial H / \partial \theta$  where

$$H = \mathbb{E} \left[ - \sum_a \pi_\theta(a|s) \ln \pi_\theta(a|s) \right] \quad (\text{S21})$$

is the Shannon entropy (note that also the alternative definition with  $\log_2$  is common). For the problem we consider, however, the choice of its coefficient  $\lambda$  is problematic: we observed that either there is no significant amplification of exploration, or the trade-off between maximizing the reward and maximizing the entropy is so strong that their common optimum is shifted away considerably from the optimum of the return. This means,  $\mathcal{R}_Q$  at the end of the simulations reaches a significantly lower value if entropy regularization is applied. We have found that it is the best strategy to start with entropy regularization for better exploration and later *smoothly* reduce its effect (by slowly decreasing  $\lambda$ ). Because we discovered this relatively close to the submission of this paper, only one of the networks shown in the main text ( fig. 7a) is trained with entropy regularization, suddenly switched off at epoch 12000.

### 3.4. Smoothening the learning gradient

We emphasize that, whenever possible, we try to consider as many cases as possible in one go, in order to reduce the noise on the learning gradient; this is like a golden thread running through the whole work presented here. Concretely, we take the following actions:

1. Our physical simulations are based on density matrices instead of wave functions. The advantage of density matrices is that they already represent the ensemble-averaged behavior of a quantum system. Dealing with wave functions would have required to sample all possible trajectories, or at least a sufficiently large subset, in order to get comparably accurate statistics. Therefore, the fact that a single density matrix contains a complete description of the quantum system prevails over the higher consumption of computational resources (density matrices are higher dimensional objects than wave functions and therefore their time evolution is numerically more expensive and they require more memory).
2. Furthermore, we also want the neural network to perform well on all possible logical qubit states. In principle, this can be achieved by training on randomly chosen logical qubit states (which requires a fair sampling and sufficient coverage of all possible states during training). Instead, however, it is much more efficient and also easier to consider the full Bloch sphere in each “trajectory”: by evolving four variables with the dimension of a density matrix ( $\hat{\rho}_0, \delta\hat{\rho}_x, \delta\hat{\rho}_y, \delta\hat{\rho}_z$  as explained in appendix 1) and rewarding the neural network based on its performance in the worst case, we directly train the network to find suitable strategies for all of them. Note that this kind of parallelism for the logical qubit state is also necessary to prevent the network from “cheating” (cmp. appendix 3.2).
3. Another point where we prefer to train rather on the average behavior than on specific examples is how measurements are treated in the reward scheme. The  $\mathcal{R}_Q$  value after a measurement can depend on the measurement result. If the agent (here, the neural network) has decided to perform a measurement, it has no further influence on the obtained result, and so it makes sense to make the reward for this action independent of this random decision. Therefore, we take the average over all possible  $\mathcal{R}_Q$  values (weighted by the corresponding probabilities to find them); appendix 3.1 describes how this can be done in a consistent way.

## 4. Recurrent network

The state-aware network develops powerful strategies for protecting the quantum information in a given qubit architecture and for a given error model. This network has, however, complete knowledge of the quantum state (or, more precisely, the completely positive map describing the evolution, represented via four density matrices, as explained in the Methods) meaning that it cannot be employed in an experiment, where information is only gained through measurements. We therefore need a neural network that is blind to the quantum state and is only aware of the measurement results. This is inherently a non-Markovian decision process: this situation requires memory to correlate previous measurements and make a decision based on their outcome. We therefore implement this state-blind network as a recurrent neural network in the form of a long short-term memory (LSTM) network [12]. The state-aware network is used as a supervisor to teach the recurrent network the correction strategy. Here we only perform supervised training but the recurrent network could in principle be trained further with reinforcement learning such that it adapts its strategy to the concrete experimental setting with a potentially uncertain knowledge about the decoherence rates or other parameters.

In contrast to feedforward neural networks with single artificial neurons as building blocks, LSTM networks consist of neurons with internal memory and more complex structure [13]. An LSTM neuron comprises a cell, input gate, output gate, and a forget gate. Each gate is a conventional artificial neuron with a nonlinear activation function that regulates the flow of information whereas the cell is responsible to remember values over time. The ability to correlate input signals over arbitrary time intervals renders LSTM networks powerful tools for processing and predicting time series in real-world applications. A more detailed description of LSTM networks can be found in Ref. [13].

The goal of our recurrent network is to learn the complete procedure of quantum error correction from the state-aware network: encoding the logical qubit, detection and recovery strategies, decoding the logical qubit, and finally, unitary operations that ensure that the final state is parallel to the initial state on the Bloch sphere. We consider the four-qubit system subject to bit-flip errors as explained in the main text. As a first step, the state-aware network is trained with immediate reward based on the recoverable quantum information and, in the last time steps, an additional immediate reward that enforces the decoding procedure is used. At the end of the time span, the network is given a reward to rotate the final state parallel to the initial state. The reward scheme is explained in appendix 3.1. After convergence is achieved, a training dataset for the recurrent network is generated.

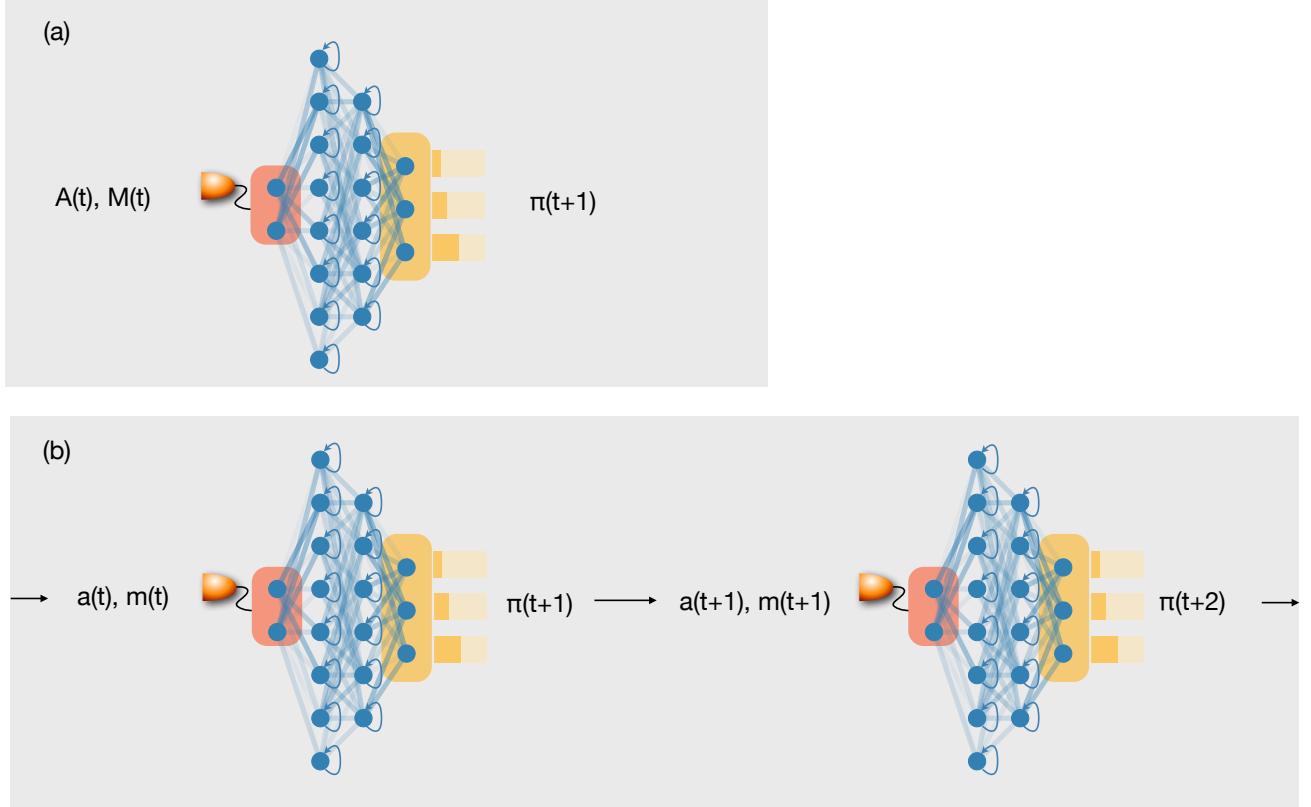


Figure S2. Visualization of the training and validation procedure of the recurrent network. (a) During training, the network is fed at the input the action  $A(t)$  and measurement result  $M(t)$  from the state-aware network, and outputs the policy  $\pi(t + 1)$  that determines the action at the next time step. This policy is trained to approximate that of the state-aware network. (b) The recurrent network is validated similarly to the state-aware network. The output policy,  $\pi(t)$ , results in the actions  $a(t)$  and results  $m(t)$ , which are then fed to the input at the next time step.

The recurrent network receives as input the action taken by the state-aware network,  $A(t)$ , and the measurement result in the  $z$ -basis (in case a measurement was performed),  $M(t)$ , and outputs the policy,  $\pi(t + 1)$ , that determines the action at the next time step, see Fig. S2(a). The goal of the recurrent network is therefore to find a non-trivial function,  $f$ , that correlates the policy at time  $t$  to the input at all previous times, i.e.,

$$\{A(0), M(0), A(1), M(1), \dots, M(t-1)\} \xrightarrow{f} \pi(t). \quad (\text{S1})$$

In the case of quantum error correction, the long-term memory is particularly important during a recovery operation after an unexpected measurement result. This is because the network needs to carefully keep track of the error syndrome and apply a correction after the recovery operation. During validation, the performance of the recurrent network is tested on the physics environment (this is the same environment that the state-aware network is trained on). In this case, its own action at time  $t$  (including the measurement outcome if the case) is fed to the input at  $t + 1$  as shown in Fig. S2(b).

The key hyperparameters of the training procedure are listed in table S1. The function that is optimized during training (the loss function) is taken to be the cross entropy between the policy of the recurrent and state-aware networks.

## 5. Hidden representation analysis via t-SNE

The goal of the state-aware network is to find an optimal encoding and correction scheme that protects the quantum information from decoherence. For bit-flip noise, the network finds various versions of the 3-qubit repetition code that are discussed at length in the main text. To understand how the strategy is reflected in the hidden representation of

Hyperparameter	Value
Architecture	Two hidden layers of 128 LSTM units
Training batch size	16
Validation batch size	256
Optimizer	Adam
Learning rate	0.001
Activation function	Hyperbolic tangent
Loss function	Categorical cross-entropy
Dropout	0.5 after each LSTM layer
Data set size	ca. 115000 gate sequences à 200 time steps

Table S1. List of hyperparameters used to train the recurrent network.

the network, we employ the t-distributed Stochastic Neighbor Embedding (t-SNE) technique [14] to the last hidden layer of the network. t-SNE is a nonlinear dimensionality-reduction technique that maps high-dimensional data to a two- (or three-) dimensional space and preserves the structure of the original data. In particular, nearby objects in the high-dimensional space are mapped to nearby points in the low-dimensional space.

In the following we give a brief description of the t-SNE algorithm. The similarity between the datapoints  $x_i$  and  $x_j$  in the high-dimensional space is modelled by a Gaussian joint probability distribution,  $p_{ij}$ ,

$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)}, \quad (\text{S1})$$

and a Student t-distribution,  $q_{ij}$ , in the low-dimensional space

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}. \quad (\text{S2})$$

The t-SNE algorithm aims at finding a low-dimensional representation of the data that minimizes the mismatch between the two probability distributions. This mismatch is quantified by the Kullback-Leibler divergence,  $C = \sum_{i,j} p_{ij} \log[p_{ij}/q_{ij}]$ , which the algorithm minimizes using gradient descent.

A critical parameter of the t-SNE technique is the perplexity,  $p$ , which is defined as

$$p = 2^{-\sum_j p_{ij} \log_2 p_{ij}}. \quad (\text{S3})$$

The variance,  $\sigma^2$ , increases as the entropy (and hence the perplexity) increases. The perplexity can therefore be interpreted as a measure of the number of neighbors considered by the algorithm. With a large perplexity, the algorithm is inclined to retain the global structure of the data, whereas a low perplexity emphasizes the local structure.

## 6. Implementation

Our RL implementation is based on the Theano framework [15], and we usually run it on GPU. Since Theano provides efficient solutions for all our neural-network related needs, we do not expect that too much can be achieved here via software optimizations. The easiest way for us to accelerate the learning part is to use one of the currently available, more powerful GPUs (currently, we have Nvidia Quadro P5000).

For the physical time evolution, we needed to implement our own simulation tools because there is no ready-to-go package meeting our specific requirements. Our current physics code is NumPy-based and runs on CPU; this is inefficient as it involves a lot of overhead and causes a memory bottleneck between physics (CPU) and network training (GPU). In addition, evolving quantum systems is a well-suited task for execution on GPU, and could therefore benefit greatly from its computing power. To improve all this, we plan a GPU-based implementation also for the physical time evolution in the future.

The numerically relevant subproblems in the time evolution scheme are (*i*) updating  $(\hat{\rho}_0, \delta\hat{\rho}_x, \delta\hat{\rho}_y, \delta\hat{\rho}_z)$  for the time evolution, (*ii*) extracting the actual network input, and (*iii*) computing  $\mathcal{R}_Q$  for the reward. Practically, (*i*) means matrix-vector multiplications between a superoperator (the completely positive map) and  $\hat{\rho}_0, \delta\hat{\rho}_x, \delta\hat{\rho}_y, \delta\hat{\rho}_z$  (interpreted as vectors); due to the sparsity of these superoperators, this operation is still relatively runtime-efficient (despite the high dimensions of the operands), and the completely positive maps describing one time step for each action of the network can be stored together in memory. For (*ii*), a principal component analysis has to be performed, and (*iii*)

	generic (non-CHZ)	CHZ-specific optimizations
time evolution	$\mathcal{O}(4^N)$	$\mathcal{O}(4^N)$
input pre-processing	$\mathcal{O}(8^N)$	$\mathcal{O}(P(N) \cdot 2^N)$
compute $\mathcal{R}_Q$	$\mathcal{O}(8^N)$	$\mathcal{O}(2^N)$
ANN training		$\mathcal{O}(P(N) \cdot 2^N)$
overall	$\mathcal{O}(8^N)$	$\mathcal{O}(4^N)$

Table S2. Runtime behavior for one epoch as a function of the qubit number  $N$ , separately for the different subproblems in our learning scheme (cmp. appendix 6).  $P(N)$  describes the dependence of PCA components on the qubit number (typically a polynomial). The training effort is assumed to be proportional to the number of input neurons.

involves the computation of trace distances. Thus, several matrices with the dimension of  $\hat{\rho}$  have to be diagonalized in (ii) and (iii).

For the special case of CHZ circuits (all unitary operations are combinations of CNOT, Hadamard and phase gates; measurement of variables in the Pauli group; state preparation in the computational basis [1, p. 464]), we can exploit the special structure of the multi-qubit quantum states that are generated. Because we still require the full density matrix, we cannot reach polynomial runtime behavior (in terms of the qubit number) like for the evolution of pure states according to the Gottesman-Knill theorem [16]; however, we can still gain an exponential factor compared to the general (non-CHZ) case. The key idea behind the efficient evolution of pure states [17, 18] is to characterize such a state as the simultaneous eigenstate of several commuting stabilizers (with eigenvalue +1 for each of them). In our adaption, we keep track of the matrix diagonal and the stabilizers determining the eigenbasis (here: stabilizers are those which commute with the density matrix). Then, (i) involves dense matrix-vector and sparse matrix-matrix multiplications (both for reduced dimensionality), and cheap update operations for the stabilizers. The effort for the diagonalization in (ii) and (iii) is completely eliminated since the matrices are already in diagonal form, and only significantly faster operations remain.

In our typical 4-qubit bit-flip examples, a reasonable level of convergence for the state-aware network is achieved within 25000 epochs. With our current implementation, the overall runtime (for the physical time evolution and the ANN training together) is around 20 hours, and the CHZ-specific optimizations reduce this time to below 6 hours. Both values are measured on the same machine (CPU: Intel Xeon E5-1630 v4, GPU: Nvidia Quadro P5000). For a well-optimized software implementation in combination with the most powerful hardware that is currently available on market, we estimate that a relative speedup of factor 10...100 is feasible.

## 7. Additional information on figures

### 7.1. Training progress due to the protection reward

In order to make the training progress in fig. 3a nicely visible, we needed to rescale the figure axes. The  $x$  axis is divided into three parts (to epoch 500, from epoch 500 to 2500, and from epoch 2500); inside these segments, the epoch number grows linearly. The data points are averaged over the last 10 epochs in the first segment, over the last 50 epochs in the second segment, and over the last 250 epochs in the third segments. The  $y$  axis is scaled as  $\mathcal{R}_Q \leftrightarrow -\ln(1 - \mathcal{R}_Q)$ , i.e. we essentially plot the negative logarithmic error probability (as  $\frac{1}{2}(1 + \mathcal{R}_Q)$  is the success rate for the optimal recovery sequence). The logarithmic scaling of the  $y$  axis is performed after taking the average.

From the  $\mathcal{R}_Q$  value at the end of a simulation, we can extrapolate an effective decay time (cmp. appendix 7.4). Fig. S3 replots the learning curve from fig. 3a in maintext, where the ticks on the right  $y$  axis indicate the effective decay time  $T_{\text{eff}}$ .

### 7.2. Discussion of trajectories of the converged state-aware network

In fig. 3(b) and (d) we show examples of gate sequences produced by the neural network. The measurement symbols indicate measurements in the  $z$ -basis and the outcome encoded in the color, i.e. dark red represents “0” or “up” and orange represents “1” or “down”. The light blue background of the quantum circuits indicates which qubits carry information about the logical qubit state. It is calculated from the criterion (Eq. (S5)). In Figure 2(c) of the main text we show the recoverable quantum information  $\mathcal{R}_Q$  as a function of time at different training stages. The results of the converged network (green lines) show a barely visible decay of  $\mathcal{R}_Q$  with time. However, a close look at this trajectory reveals interesting substructure. Therefore, in Fig. S4(a) we show this trajectory (labeled “traj. 1”) from

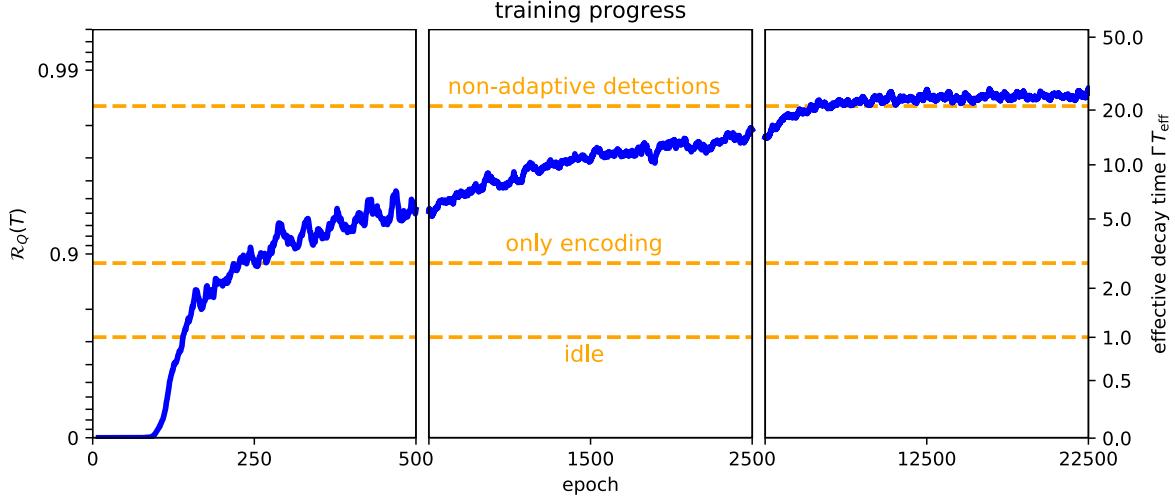


Figure S3. Replot of fig. 3a with a second  $y$  axis indicating the effective decay time  $T_{\text{eff}}$  extrapolated from the decay of  $\mathcal{R}_Q$  during the simulation time ( $T = 200\Delta t$ ). Note that this is a well-defined quantity only for strategies where  $\mathcal{R}_Q$  decays exponentially (or at least with an exponential envelope) which is in particular not the case for “only encoding”.

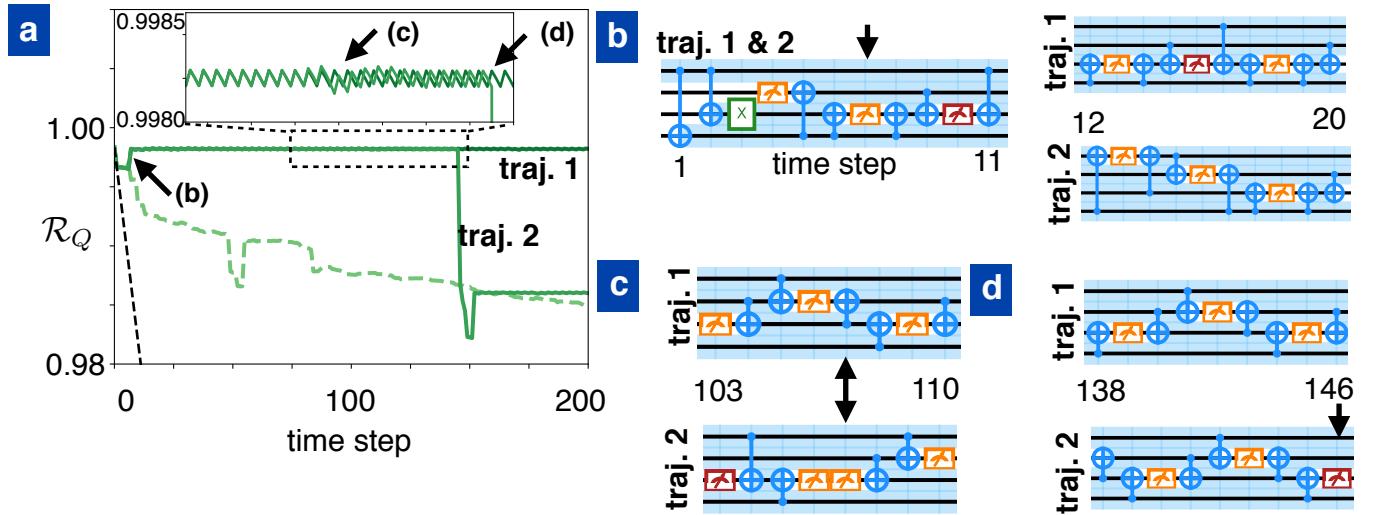


Figure S4. (a) The recoverable quantum information  $\mathcal{R}_Q$  as a function of time simulated with a converged neural network. Two typical trajectories (solid lines) and the average over a batch of 128 trajectories (dashed) are shown. The trajectory labeled “traj. 1” is the same as shown in fig. 3c of the main text. (b) The gate sequence (similar for both traj. 1 and traj. 2) responsible for the initial drop of  $\mathcal{R}_Q$ . Only after the successful encoding (time step 2) the decay is slowed down. In time step 7 the first parity measurement leads to an increase of  $\mathcal{R}_Q$ . After time step 11 the two trajectories show different gate sequences that, however, are equivalent in terms of preserving  $\mathcal{R}_Q$ . (c) The gate sequence around the time where a shift between the respective  $\mathcal{R}_Q$  is observed. In trajectory 2 the network repeats one measurement before continuing with repetitive parity measurements, while the other trajectory sticks to its original periodicity. (d) The gate sequences where trajectory 2 experiences an unexpected measurement result (indicated by an arrow). This leads to a drop in  $\mathcal{R}_Q$  which is partially recovered after a few further parity measurements that pinpoint the error to one specific qubit.

the main text, the corresponding sample-average, as well as another typical trajectory. Here it becomes visible, that both trajectories are subject to a fast decay in the first few time steps, where the encoding does not yet slow down the normal decoherence process. After distributing the information about the logical qubit on several physical qubits, the effective decoherence rate is reduced, but only after the first parity measurement  $\mathcal{R}_Q$  increases again. Note that this increase of the recoverable quantum information is not guaranteed, but depends on the measurement outcome: If the measurement result indicates that the two qubits are still in the same state (which is the most likely outcome after a sufficiently short time span) then the one can safely assume that no error occurred so far and a small revival of  $\mathcal{R}_Q$  is observed. However, if by chance the measurement outcome indicates a bit flip, a drop of  $\mathcal{R}_Q$  would occur instead. The corresponding gate sequences are displayed in Fig. S4(b). The network immediately performs the CNOTs to switch to the encoded state, where the decay of  $\mathcal{R}_Q$  is slowed down. Curiously, this particular network consistently performs two unconventional additional actions (a bit flip and a measurement), before preparing the first parity measurement. (training another neural network does not reproduce these two actions). One could speculate, that the seemingly unnecessary very first measurement of the second qubit is performed in order assure the state of this qubit since it is then used as the measurement ancilla.

In the inset of Fig. S4(a) it becomes visible, that the subsequent repeated parity measurements lead to a periodic tiny decay and recovery of the recoverable quantum information, since in these trajectories, all the measurement outcomes indicate that no error occurred. Note that, although a different sequence of parity measurements is performed for trajectory 1 and 2, the recoverable quantum information in both cases is the same. The slight shift of the trajectories that eventually appears is due to a repeated measurement in trajectory 2, see Fig. S4(c). Only when an error is detected, as in trajectory 2 (Fig. S4(d)), the behavior changes: The recoverable quantum information drops, until additional parity measurements determine which qubit was actually subject to the bit-flip. Then  $\mathcal{R}_Q$  partially recovers. Note that these are tiny small-scale drops compared to the large-scale drops visible for the not yet converged network in fig. 3c of the main text (which are cause by more than single bit-flip errors or unclear situations because too much time passed between consecutive measurements).

The dashed line in Fig. S4(a) shows the sample average over 128 trajectories. It decays faster than trajectory 1, where no error occurred at all due to all the trajectories with bit flips (and a corresponding drop in  $\mathcal{R}_Q$ ) at various different times that enter the average (note that the “revival” of the sample average is only an artefact due to the small number of trajectories used here).

### 7.3. t-SNE

In fig. 4 and fig. 7 of the main text, we show graphical representations of the quantum states via the yellow rectangles. The border of the rectangles is colored according to the action the network takes for the particular input. In detail, we show the six states with the largest contribution to the density matrix  $\hat{\rho}_x$  ordered by their associated probability. The exact value of the probability is visualized by the blue bars next to the states. The states with a probability less than 5% are semi-transparent. This represents as a showcase how an initial superposition state  $(|0\rangle + |1\rangle)/\sqrt{2}$  evolves. The black (white) circles correspond to 1 (0) so that, e.g., the state  $\bullet \bullet \circ \circ + \bullet \circ \circ \bullet$  is equivalent to  $(|0100\rangle + |1001\rangle)/\sqrt{2}$ .

Here we apply the t-SNE technique to the hidden representation of a fully converged neural network that is trained via reinforcement learning to protect four physical qubits against bit-flip noise as explained in the main text. The network has three hidden layers of 300 neurons each. To improve speed, we first perform a principal component analysis that maps the original 300-dimensional data onto a 50-dimensional space. Aside from perplexity, all other parameters of t-SNE are kept at their default value [19]. To highlight both the global and the local structure of the hidden representation, we apply the t-SNE algorithm using three different values of perplexity (compared to the data set): small, intermediate, and large. In the following we discuss the map presented in the main text, which is computed for an intermediate value of perplexity and captures a subtle interplay between global and local features. Later we will discuss two additional maps with large and small perplexity.

*Intermediate perplexity.* In fig. 4 in the main text, we perform the t-SNE algorithm on a network validated on 1024 trajectories, each with 200 time steps, which yields a data set of 204,800 points. A perplexity of 2500 is used. A point on the map corresponds to the 2D-projected hidden representation for a particular input (the four density matrices of the RL-environment), and is colored according to the action the network takes. The map features several dense clusters as well as smaller islands and collection of points. To understand this pattern, we study a steady-state detection sequence that is visualized in fig. 4a. An inspection of the density matrices reveals that, in steady state, one quantum state has the largest contribution, which is the state the network aims to protect from decoherence. The other states are the leading error channels that the network aims to remove by applying the steady-state detection cycle: CNOT(1,3)CNOT(4,3)M(3)CNOT(1,3)CNOT(2,3)M(3). In this case, the network uses the second qubit as the

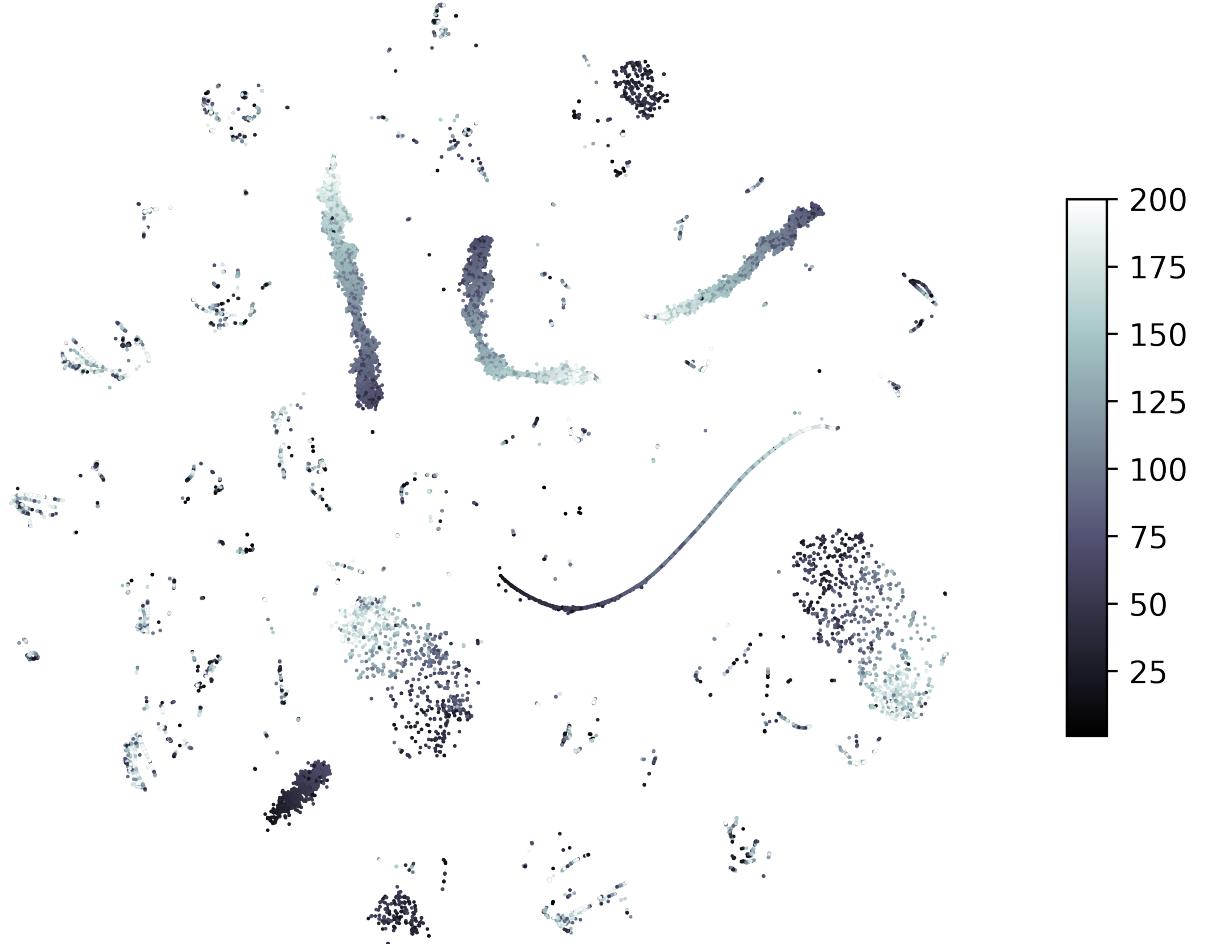


Figure S5. Same t-SNE map as in fig. 4a in the main text but colored according to time. The arrow of time is clearly discernible in the dominant clusters owing to the slow decoherence of the density matrix.

ancilla and an encoding of the quantum information in the other three qubits of the form

$$|\Psi_{\text{enc}}\rangle = \alpha|011\rangle + \beta|100\rangle \quad (\text{S1})$$

The density matrix can be uniquely identified at any time during the steady-state detection cycle. For instance: (i) the states after a measurement have a higher coherence than before the measurement, (ii) after a cnot operation the ancilla is entangled with the other qubits, and (iii) the orange and purple nuclei correspond to a different state of the ancilla qubit. These differences are clearly understood and exploited by the neural network judging from the clusters in the hidden representation.

The detection strategy is six-fold periodic and can be visualized in the hidden representation via the gray trajectory in fig. 4c. The nodes of the trajectory are displayed as white circles. The trajectory features jumps between the six largest clusters of the map, which contain most of the points of the map. This is because unexpected measurements (which interrupt the detection cycle and take the network outside the main clusters as described below) are improbable events: using the parameters from the simulation we obtain a probability of  $\sim 50\%$  over the entire trajectory that all measurements on the ancilla qubit yield the expected result. A peculiar feature of the main clusters is their strongly asymmetric shape, which is due to the slow drift of the density matrix subject to decoherence. This can be visualized in the inset of fig. 4c, where the orange cluster on the right-hand side of the map is colored according to a function that

varies linearly with time (white corresponds to  $t = 1$  and dark orange to  $t = 200$ ). The detection strategy therefore forms nearly closed loops in the hidden representation that slowly drift as time progresses. Some clusters are spread out so much that they are “cut” by other clusters. An example can be seen in Fig. S5, where the map is colored according to time. Comparing this map with the map in the main text one notices the dominant CNOT13 cluster on the left-hand side of the map being cut by the M3 cluster. This “artifact” is likely due to the fact that the perplexity is not sufficiently high and the global structure is not fully captured.

The quantum system is brought out of the steady state by an unexpected measurement result, which means that a measurement on the ancilla qubit does not project the quantum system onto the state with the largest contribution (see the discussion in the above paragraph). This interrupts the periodic detection sequence; the network employs a fully adaptive strategy in this regime that is carefully optimized during training since it plays a key role in the capability to recover the quantum information. This recovery strategy deals with markedly different quantum states and thus corresponds to different trajectories in the hidden representation. An example of a recovery trajectory is displayed by the blue line in fig. 4c, where the white stars denote the nodes of the trajectory. These points belong to small clusters outside the dominant clusters.

In the following we discuss key features of this particular recovery strategy. At step 1, an unexpected measurement result projects the quantum state onto a mixture of three pure states with comparable contribution, which correspond to a bit flip of qubit 2 (largest contribution), qubit 1 (second largest), and qubit 3 (i.e., the ancilla qubit), see fig. 4b. The probability that qubit 2 flipped is roughly twice as large as the one of qubit 1 and 3 (whose probability is about the same). This can be understood as follows. Consider the state at the beginning of the detection cycle (e.g., belonging to the blue cluster in the figure)

$$|\Psi\rangle = \alpha|011_A1\rangle + \beta|101_A0\rangle, \quad (\text{S2})$$

where the subscript “A” denotes the ancilla. The ancilla is the physical qubit that has been designated by the network to be measured. The first part of the detection cycle, CNOT(1,3)CNOT(4,3)M(3), projects the ancilla state onto  $|0_A\rangle$  assuming an expected measurement result. The error channels in which qubits 1, 2, or 4 flip map the ancilla state into  $|1_A\rangle$ ,  $|0_A\rangle$ ,  $|1_A\rangle$ , respectively, meaning that an expected measurement result  $|0_A\rangle$  leaves the flip of qubit 2 to be the leading error channel. The second part of the detection cycle, CNOT(1,3)CNOT(2,3)M(3), maps the ancilla state back onto  $|1_A\rangle$ , while the error channels are mapped onto  $|0_A\rangle$ ,  $|0_A\rangle$ , and  $|1_A\rangle$ , respectively. Clearly, in the case of an unexpected measurement result  $|0_A\rangle$ , the three major contributions are the flips of qubits 1, 2, and the ancilla itself. Since only qubit 2 was not corrected after the first part of the detection cycle, it carries a contribution that is twice as large. These features can be clearly seen in the density matrix.

After projection onto an unexpected state (step 1), the network begins a recovery cycle that should ideally be as short as possible. The network enforces this by targeting the state with the largest contribution in the density matrix. As a result, application of CNOT(4,3)CNOT(2,3)M(3) (steps 2–4) intentionally swaps the ancilla of the leading contribution to  $|1_A\rangle$ . At this particular time step, the outcome of the measurement yields  $|0_A\rangle$  meaning that the network needs to perform another cycle (steps 5–7) to bring the quantum system back to a coherent state. After the measurement at step 7, it is finally revealed that the ancilla was flipped.

*Large perplexity.* In the following we perform the t-SNE algorithm on a network validated on 16 trajectories, each with 200 time steps, which yields a data set of 3200 points. To faithfully reproduce the global structure of the representation, we use a relatively large perplexity (compared to the data size) of 200. Upon a visual inspection of the map in Fig. S6, six main clusters can be identified that likely reflect the six-fold periodicity in the detection scheme. Note that after the next layer (i.e., the output layer), these clusters are converted to nearly deterministic actions. Each cluster has a non-trivial substructure with a nucleus of tightly packed points and peripheral domains with smaller density. We identify the nuclei with the steady state of the quantum system in which measurements yield the expected outcome and increase the coherence. This is the most abundant situation due to the small decoherence rate used in the simulations and, thus, a small probability that a measurement yields a result with an unexpected outcome. Consequently, the nuclei contain the majority of the points of the cluster.

Each nucleus has a representative input (density matrix) that is shown in Fig. S6(a). A steady-state detection strategy is six-fold periodic and is reflected in the hidden representation as jumps between the nuclei. An inspection of the density matrices in fig. 4b reveals that, in steady state, one quantum state has the largest contribution, which is the state the network aims to protect from decoherence. The other states are the leading error channels that the network aims to remove by applying the steady-state detection cycle, see the discussion from the intermediate-perplexity section.

Despite being tightly packed, the nuclei show some substructure as well due to the fact that the quantum system decoheres over time. This slow drift of the density matrix can be seen in Fig. S6(c), where the nucleus belonging to the red cluster is colored with a function that depends linearly on time, from  $t = 1$  (white) to  $t = 200$  (black).

The peripheral domains of a cluster belong to states in which the quantum system is brought out of the steady state by an unexpected measurement result, which means that a measurement on the ancilla qubit does not project

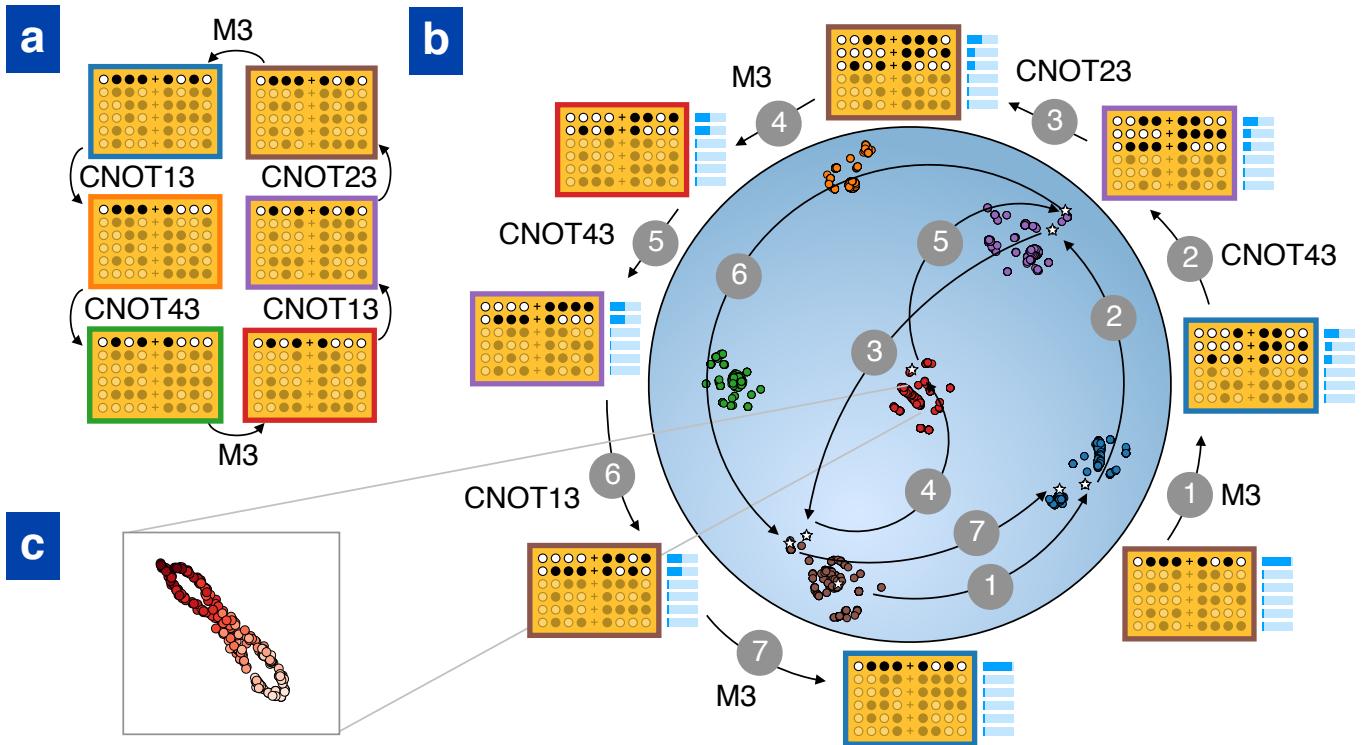


Figure S6. Visualization of the hidden representation of the neural network using the t-SNE technique in the large-perplexity regime. (a) Sequence of states visited during a standard repetitive detection cycle. (b) Visualization of neuron activations (300 neurons in the last hidden layer), sampled in several runs, projected down to 2D using the t-SNE technique. For a particular gate sequence that is triggered upon encountering unexpected measurements, we also indicate the states and the actions producing transitions between the states. Qubits are numbered 1,2,3,4, and CNOT13 has control-qubit 1 and target 3. (c) Zoom-in shows a set of states in a single cluster; the shading indicates the time progressing during the gate sequence, with a slow drift of the state due to decoherence.

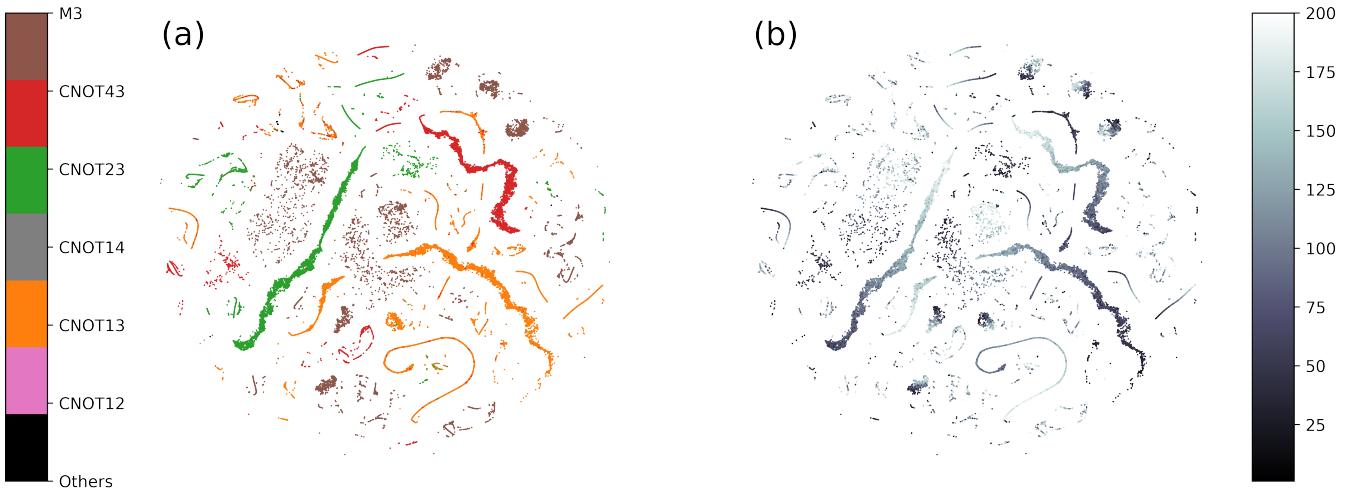


Figure S7. Visualizing the hidden representation of the state-aware neural network with the t-SNE technique in the small-perplexity regime. In (a), the map is colored according to the action that the neural network is about to take. The large circles denote the points for which the coherence of the state is low; this happens after an unexpected measurement result occurs and lasts for three or six time steps depending on the results of the next measurements, see text for details. On the right, the coloring is performed according to the time. The long and thin clusters correspond to the (slow) drift of the density matrix over time.

the quantum system onto the state with the largest contribution. The recovery cycle can also be visualized in the hidden representation of the network, where the white stars denote the points belonging to the cycle. Initially (before step 1) the quantum system is in steady state and the point belongs to the nucleus of the brown cluster. The points during the recovery operation belong to the peripheral domains. The corresponding density matrices have common features with the cluster they belong to, yet are sufficiently dissimilar and are clearly separated from the nuclei. After the recovery cycle (after step 7), the quantum system goes slowly back to steady state and no longer ends up in a nucleus but rather in another conglomerate of points close to it. We have observed this to be a generic feature of all trajectories and is caused by the slight decrease of coherence of the states after a recovery cycle, which the network observes through the density matrix.

*Small perplexity.* In the following we use the t-SNE technique to visualize the hidden representation in the small-perplexity regime for getting a deeper insight into the local structure of the data. We evaluate 1024 trajectories of 200 time steps each yielding a dataset of 204,800 points with a perplexity of 400. The resulting map is displayed in Fig. S7 and features several long and thin clusters that would form the nuclei in the high-perplexity regime as well as smaller and less dense clusters that would likely form the peripheral domains. In Fig. S7(a) the map is colored upon the action the network is about to take. There are six actions the network takes: four are related to the detection and recovery strategies mentioned above, and the other two, CNOT(1,2) and CNOT(1,4), are performed at the very beginning of the trajectory to encode the logical qubit. The slow decay of the coherence is revealed particularly well on this map, see Fig. S7(b). Despite revealing the local structure of the map with great detail, the global information is somewhat distorted. For instance, note the long orange stripe on the right-hand side of the map in Fig. S7(a) that is “cut” towards the middle by another cluster. This cut becomes particularly obvious if one follows this cluster in Fig. S7(b).

#### 7.4. Effective decoherence time

In fig. 5(a) and (b) of the main text we show different qubit connectivities that are realized by restricting the set of available gates that the neural network can use. Only single qubit gates and CNOTs between connected qubits can be performed. We also restrict the qubits on which measurements are available. In the gate sequences shown in fig. 5a the effect of a geometrical restriction becomes clearly visible: A simple parity measurement requires now (on average) more than three operations (which would typically be two CNOT and one measurement). For example, in the case of the circular geometry, where measurements are only allowed on an ancilla qubit with a single nearest neighbour it requires at least six operations to entangle two different qubits to the ancilla. One of those can be the nearest neighbour of the ancilla, which requires a single CNOT to entangle it, but the second qubit state needs first to be swapped with the nearest neighbour state (three CNOTS), followed by the entanglement (one CNOT) and the actual measurement. Similarly also for the chains, the gate sequence for the repetition code becomes effectively longer than compared to the all-to-all connected system. This increases the average time between consecutive parity measurements which affects both the slow-down of the overall decay of  $\mathcal{R}_Q$  as well as the recovery of  $\mathcal{R}_Q$  after an error actually occurred. The physical consequences of preserving the recoverable quantum information as good as possible can be understood in terms of a prolonged effective decoherence time  $T_{\text{eff}}$ . Therefore, depending on the average length of the gate sequences necessary to perform parity measurements and pinpoint errors, we expect different effective decoherence times for the different configurations. In fig. 5b of the main text, we show this effect. In particular, we extract the effective decoherence time  $T_{\text{eff}}$  from the average decrease of  $\mathcal{R}_Q$  after the total time  $T$  and plot the ratio  $T_{\text{eff}}/T_{\text{dec}}$  where  $T_{\text{dec}}/\Delta t = 1200$ . We use that the recoverable quantum information decays approximately exponentially with time, i.e.,

$$\langle \mathcal{R}_Q(T) \rangle = \exp \left( \frac{-2T}{T_{\text{eff}}} \right), \quad (\text{S3})$$

where  $\langle \mathcal{R}_Q(T) \rangle$  denotes the sample average (over 10 000 trajectories) of the recoverable quantum information at the final time  $T$ . The result show a clear trend: While all configurations increase the effective decoherence time and thus provide additional stability against bit-flip errors, the all-to-all connected systems allows for the largest improvement, followed by the chain where the neural network can use all qubits for measurements. Notably, this chain configuration performs better than the chain with one measurement qubit only (located on a central spot, at the edge would be even worse). This is, because in the first case the network switches between several qubits for measurements which leads to an advantage as compared to the fixed location. The circular configuration clearly requires the longest gate sequences and consequently it achieves a smaller increase of  $T_{\text{eff}}$  than the other systems.

### 7.5. Introducing measurement errors

In fig. 5c of the main text, we give a further example how neural networks can adapt their strategies to different environments which they are trained on. In particular, while in all other figures we have used perfectly reliable measurement outcomes, we now introduce a measurement error probability and see how a network trained with this new environment learns an adapted strategy instead of converging to the same solution as before.

We implement measurement errors by modifying the form of  $\phi$  (cmp. Eq. (S4)) for measurements: for completely reliable measurement results as considered in all other simulations, measurement result  $j$  is associated with  $\phi_j(t + \Delta t, t)[\hat{\rho}] = e^{\Delta t \mathcal{D}}(\hat{P}_j \hat{\rho} \hat{P}_j^\dagger)$  (including the dissipative dynamics during the following idle time  $\Delta t$ , cmp. Methods), and for measurement errors this expression is modified to

$$\phi_j(t + \Delta t, t)[\hat{\rho}] = e^{\Delta t \mathcal{D}} \left( \left( \sum_k p(R_j|A_k) \hat{P}_k \right) \hat{\rho} \left( \sum_k p(R_j|A_k) \hat{P}_k \right)^\dagger \right) \quad (\text{S4})$$

where  $p(R|A)$  is the probability that result  $R$  is reported given that the quantum system is actually projected into state  $A$ . Explicitly, we considered ‘‘mixing matrices’’ of the form

$$\begin{pmatrix} p(0|0) & p(0|1) \\ p(1|0) & p(1|1) \end{pmatrix} = \begin{pmatrix} 1 - \epsilon & \epsilon \\ \epsilon & 1 - \epsilon \end{pmatrix} \quad (\text{S5})$$

for different values of the measurement error probability  $\epsilon$  between 0.001 and 0.1. For each value of  $\epsilon$  shown in fig. 5c we trained a separate neural network from scratch.

We find that the neural network counters increasing measurement errors by repeating the same measurement for several times. This helps to identify the true measurement outcome and to avoid wrong conclusions. On the other hand, repeating a measurement increases the time until the next parity measurement where new information about errors can be gained. Thus, depending on the percentage of false measurements, the neural networks learn to perform only a few or many additional measurements. We show the fraction of measurements, i.e. how many actions of all actions are measurements, in fig. 5c of the main text. This data is obtained as an average of  $2^{13} = 8192$  samples and we also indicate the standard deviation (for the first data point it is so small that the error bars vanish behind the dot). We observe that the average number of measurements used increases significantly. Here, in Fig. S8 we show a typical trajectory for each measurement error probability that was used. In the lowest two trajectories, where the error probability is already very large, one can readily identify examples of repeated measurements where some show a differing, false outcome. The anomalous outcome is specific to one measurement in contrast to a bit-flip that would also affect all future measurements.

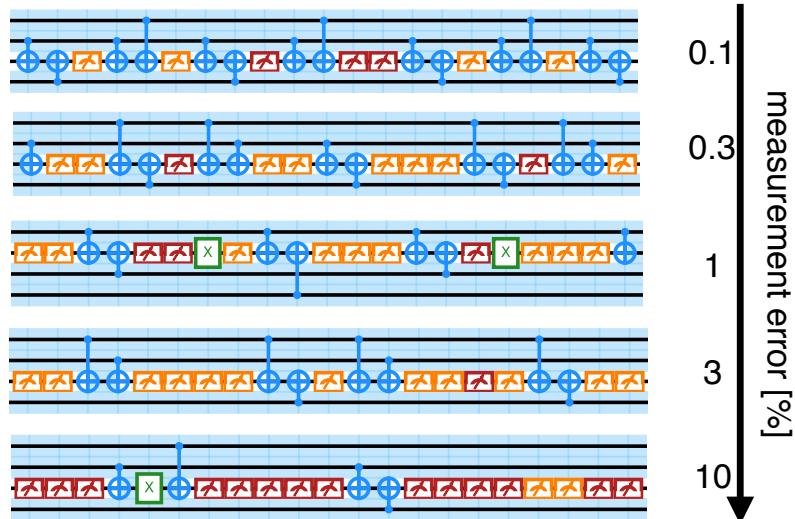


Figure S8. Typical gate sequences obtained from networks that were trained with different measurement error probabilities.

### 7.6. The behavior of the state-aware network with recovery reward

According to the recovery reward structure described in appendix 3.1, only errors of the specified target qubit need to be corrected to earn the correction reward. Errors occurring on any other qubit might or might not be corrected, as long as their status of being flipped or not is remembered for the interpretation of further parity measurements and the final decoding. Since the correction reward is only obtained after all time steps, in principle the network could decide to correct also the errors of the target qubit only at the very end (after decoding). However, we observe that the network corrects errors on the target qubit after an intermediate time - it typically performs a few further parity measurement to determine the qubit where the error occurred. Even when this is clarified with more than 90% certainty, it usually performs a few more parity measurements. We believe that this is done to optimize the immediate reward earned due to preserving the recoverable quantum information. This is more pressing (since otherwise reward might be lost) than the actual correction which will only be rewarded in the very end. Furthermore, we observe that sometimes also errors on other qubits are corrected. This might be related to our observation of some “corrections” even in the absence of a correction reward. There we conjecture that the network performs these corrections in order to return to a well-known state.

While the correction is a rather slowly learned property (cf. Fig. S9), decoding is learned rapidly and with very close to perfect success. We observe that the network immediately starts decoding after receiving the signal and fills up all remaining time steps with actions that are equivalent to idle, since they do not operate on the qubit that carries the logical qubit state (e.g. measurements on already decoded qubits, cf. inset of fig. 7a in the main text where the error indicates the time step where the decoding signal is started). This is due to our reward structure that punishes the network if it entangles additional qubits to continue error correcting. In practice, of course, one would want to immediately read-out or proceed to perform operations on the decoded qubit, since it no longer profits from any protection.

In Fig. S9(a) we show the mean of the rewards earned during the first 2000 training epochs. We show the reward split up in its four components and averaged over the time  $T$  and a batch of 128 different trajectories: the mean of all immediate rewards earned at individual time steps for preserving the recoverable quantum information (blue); the mean reward (negative) for performing destructive measurements (black); the mean of all rewards earned for performing decoding (green); the mean correction reward (black). The shaded areas indicate the corresponding standard deviation which gives an approximate understanding of how strongly the individual parts contribute to the learning gradient.

In Fig. S9(b) we show how successfully the different aspects of the full quantum error correction were learned: In blue we show the recoverable quantum information  $\mathcal{R}_Q$  after the last time step averaged over the validation batch. In black we show the fraction of trajectories in the validation batch where no destructive measurement (destroying the logical qubit state) was performed. In green we show the fraction of trajectories where the decoding into the desired target state was successful (checked at the last time step). And in red we show a measure of the overlap  $\mathcal{O}_Q$  to indicate how well the initial logical qubit state has been preserved overall (for the definition, see Methods).

It is interesting to note that the overlap criterion rises much slower than all the other quantities. In particular, the network learns rather early and very fast to avoid destructive measurements and to decode properly. This is achieved with a very close to perfect success rate. There is a fast increase in preserving the recoverable quantum information at the very beginning which is associated to avoiding destructive measurements and finding a good encoding. This is followed by a slow convergence (not shown in the figure but for much longer training times of several 10 000 epochs) where the network optimizes the error correcting sequence by avoiding non-destructive but also not helpful actions. Also the overlap measure converges much slower, indicating that correcting actual bit flips consequently is learned only at later training times. This is partially related to the fact that it can only be learned if the recoverable quantum information is already decently preserved and the decoding is performed properly, as well as to the fact that important bit flips (on the target qubit for the logical qubit state) that really need to be corrected occur not that often. The decoding is very early on learned with a high success rate, i.e., typically all trajectories of a batch (with very few exceptions that become even rarer during later training stages) end up correctly decoded. The large standard deviation indicated around the decoding reward despite its already high overall success is due to the fact that the decoding steps are not necessarily performed at the same time steps for all trajectories (even if at the end of the time all trajectories end up in the correct state).

### 7.7. Validating the recovery success

We need a suitable measure to judge the success of the recovery sequence. In our (bit-flip) applications, the remains of the logical qubit state (after successful decoding) should be stored only in one target qubit with maximum overlap to the (pure) target state  $|\phi_{\vec{n}}\rangle$ . This is quantified by  $\langle \phi_{\vec{n}} | \text{tr}_{\text{o.q.}}(\hat{\rho}_{\vec{n}}) | \phi_{\vec{n}} \rangle$  ( $\text{tr}_{\text{o.q.}}$  denotes the partial trace over all qubits

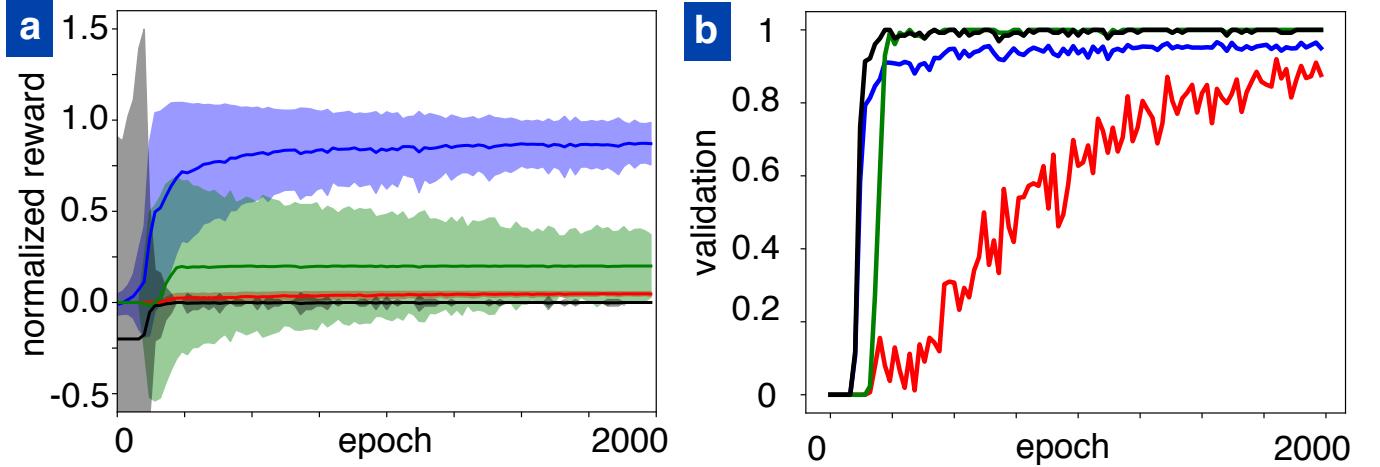


Figure S9. (a) The normalized rewards earned on the validation batch as a function of training epochs. Blue: The sum of all immediate rewards earned for preserving the recoverable quantum information. Black: The average (negative) reward obtained due to performing destructive measurements. Green: The average total decoding reward. Red: The average correction reward earned at the last time step. The shaded area indicates the standard deviation, where the observed roughness (especially the vanishing and reappearing black shading) is due to the relatively small batch size of 128 trajectories. (b) The training success validated on a batch of 128 samples as a function of training epochs. Blue: The average recoverable quantum information  $\mathcal{R}_Q(T)$  after the last time step. Black: The fraction of trajectories where no destructive measurement was performed. Green: The fraction of trajectories, where the decoding to the desired target state was successful (after the last time step). Red: As measure for the overlap,  $2(\mathcal{O}_Q - 1/2)$  that indicates how well the logical qubit state is reproduced after the last time step.

except for the target qubit), and to describe the worst-case scenario over the whole Bloch sphere, we consider

$$\mathcal{O}_Q = \min_{\vec{n}} \langle \phi_{\vec{n}} | \text{tr}_{\text{o.q.}}(\hat{\rho}_{\vec{n}}) | \phi_{\vec{n}} \rangle \quad (\text{S6})$$

(cmp. Eq. (S7)). This defines how well the final state preserves coherence *and* how well it matches the target state in terms of direction. Since ‘‘random guessing’’ already leads to a value of  $\mathcal{O}_Q = \frac{1}{2}$ , we rather plot  $2(\mathcal{O}_Q - \frac{1}{2})$  in fig. 7(a) and (b) of the main text.

### 7.8. Performance of the recurrent network

In fig. 7b of the main text we plot the relevant validation quantities for the recurrent network during training. Those are the recoverable quantum information at the final time,  $\mathcal{R}_Q(T)$ , and a measure of the overlap,  $\mathcal{O}_Q$ . Both quantities rise monotonically during training and converge towards the level of the state-aware network. Interestingly, preserving the quantum information (captured by  $\mathcal{R}_Q$ ) is learned faster than rotating the logical qubit (captured by  $\mathcal{O}_Q$ ). This is likely because no long-term memory is required to preserve  $\mathcal{R}_Q$  (the periodicity of the error detection sequence is six time steps only) whereas rotating the logical qubit requires memorizing the entire recovery procedure and therefore demands longer-term memory (in fig. 7c the correction is applied thirteen time steps after the unexpected measurement result).

The key aspects of the strategy learned by the recurrent network are visualized in the quantum circuit of fig. 7c during validation. The usual detection strategy is a modified version of the repetition code in which the network alternates the measurement qubit. The encoded state of the logical qubit is of the form  $\Psi_{\text{enc}} = \alpha|000\rangle + \beta|111\rangle$ . After the unexpected measurement result visualized at the left of the black rectangle at  $t = 1$  in the figure, the network initiates the recovery procedure. As usual, mainly three states contribute to the density matrix at this time as seen in the figure: the state in which only qubit 3 flipped (largest contribution), qubits 2 and 3 flipped, and qubits 1 and 3 flipped (smallest contribution). Note that the density matrix is only plotted for convenience — the recurrent network does not have access to it. After  $t = 2$ , the network copies the quantum information from qubit 2 to 3, and after  $t = 3$  checks whether qubit 2 flipped. At  $t = 5$ , the state in which qubits 2 and 3 flipped is ruled out by measurement. Afterwards, the two remaining states are compared and, at  $t = 7$ , the measurement confirms that qubit 1 flipped. In principle, the network has sufficient information to apply a bit flip to qubit 1 at this point. However, it turns out that after this six-step recovery, the quantum state has a slightly lower coherence than in steady state. Between  $t = 7$  and

$t = 10$ , the network therefore conducts a further purification procedure. The recovery procedure ends at  $t = 13$  when a correction to qubit 1 is applied.

To understand how the strategy is reflected in the internal operation of the recurrent network, we study the hidden representation of the last hidden layer. In fig. 7c we plot the output of the selected LSTM neurons during the above-described detection and recovery procedure. The hidden representation and quantum circuit are aligned such that the neuron activations result from feeding the corresponding action to the input of the network. The first six neurons (1c–6c) are mainly involved in the steady-state detection procedure judging from their periodic activations before the unexpected measurement result. Each of these neurons fire whenever a specific gate of the detection sequence is applied. The neurons involved in the recovery process (1r–4r) display correlations over longer times:

1. Neuron 1r starts firing after an unexpected measurement outcome and keeps firing until the final correction to qubit 1 is applied. This neuron presumably tells the network whether a recovery procedure is currently underway.
2. Neuron 2r fires particularly strongly if a measurement on qubit 3 yields an unexpected result. In contrast, neuron 3c fires if the same measurement yields an *expected* result.
3. A more peculiar firing pattern can be observed for neuron 3r: during steady-state detection the neuron displays a smooth and periodic activation, which the network likely uses as an internal clock during the standard sequence. However, during recovery the pattern changes and the neuron fires when the gate CNOT(1,2) is applied. Note that this gate is only applied during recovery (it may also be applied during encoding and decoding but not during detection) and signals that qubit 2 is now disentangled and should be measured at the next time step.
4. A dual role is played by neuron 4r as well. During steady-state detection it fires if CNOT(1,4) is applied. During recovery, however, it fires whenever qubit 2 is used as ancilla. This provides relevant information to the network regarding the encoding of the logical qubit.

### 7.9. Correlated Noise Scenario

To further verify the flexibility of our approach, we consider in fig. 6 of the main text a different error model than in the rest of the result discussion: in fig. 6, we have correlated noise described by the Lindblad equation (cmp. Methods)

$$\frac{d}{dt}\hat{\rho} = \frac{1}{T_{\text{dec}}}\left(\hat{L}\hat{\rho}\hat{L}^\dagger - \frac{1}{2}\left\{\hat{L}^\dagger\hat{L}, \hat{\rho}\right\}\right) \quad (\text{S7})$$

with decoherence time  $T_{\text{dec}}$  and the Lindblad operator

$$\hat{L} = \frac{1}{\sqrt{\sum_q \mu_q^2}} \sum_q \mu_q \hat{\sigma}_z^{(q)} \quad (\text{S8})$$

where  $q = 1, 2, 3, \dots$  labels the qubit,  $\mu_q$  is the corresponding magnetic moment, and  $\hat{\sigma}_z^{(q)}$  is the Pauli  $z$  operator for the respective qubit. In this discussion, we will always assume that qubit 1 is the data qubit; the remaining ones are called the ancilla qubits.

The action set consists of measurements on these ancillas along the  $x$  and  $y$  axis, plus the idle operation; note that we do not allow CNOT gates here. The described setup allows to extract information about the noise on the data qubit from measurements on the ancillas. The agent can decide in which order these measurements (regarding both the qubit and the axis) are performed, and how many idle steps are chosen in between; a priori, it is not clear which strategy is the best. As we will see, this question becomes very complex for more than one ancilla, especially because improvements can be achieved by employing adaptive feedback schemes (further action sequence depending on previous measurement results).

For a benchmark of the behavior of our neural network, we extrapolate an effective decay time from the averaged value  $\langle \mathcal{R}_Q(T) \rangle$  at the end of the simulations, assuming an approximately exponential decay. These values are directly comparable to the theoretical results that we will compute below. To make their interpretation easier, we always normalize by the “trivial” decay time  $T_{\text{triv}}^{(\text{cn})}$  for a single qubit (see Eq. (S9)), and call this ratio  $T_{\text{eff}}/T_{\text{triv}}^{(\text{cn})}$  the (coherence) improvement.

In order to actually judge how well our networks perform on this problem, we will in the following describe alternative approaches to find suitable measurement schemes in the two- and three-qubit scenario. We emphasize that already for this problem, the effort in terms of manpower for these alternative ways exceeds considerably that for training a neural network.

*Single qubit scenario* If all ancillary qubits are ignored (mathematically described by a partial trace over all qubits except for the data qubit), or equivalently if there are no ancillas,  $\mathcal{R}_Q(t) \sim \exp(-2t/T_{\text{triv}}^{(\text{cn})})$  with the “trivial” decay time

$$T_{\text{triv}}^{(\text{cn})} = \frac{\sum_q \mu_q^2}{\mu_1^2} \cdot T_{\text{dec}} \quad (\text{S9})$$

(called  $T_{\text{single}}$  in the Methods).

*1+1 qubit scenario* With one ancilla qubit, it is already possible to achieve a slowdown in the decay. This two-qubit case can be treated fully analytically.

The best strategy is to measure the  $x$  and  $y$  axis of the ancilla qubit in an alternating manner. To get insights into the effect of the intermediate idle time  $\tau$ , we will derive a closed form for the effective decay time  $T_{\text{eff}}(\tau)$  in the following; the result, Eq. (S13), is used to plot the analytical predictions in fig. 6b of the main text.

We start by investigating the dynamics of the quantum system if this protocol is applied. Immediately after each measurement at time  $t_m$ , the quantum system (for logical qubit state  $\vec{n} = (x, y, z)$ ) is in a product state of the two qubits:

$$\hat{\rho}_{\vec{n}}(t_m) = \frac{1}{2} \begin{pmatrix} 1 + z_m & x_m - iy_m \\ x_m + iy_m & 1 - z_m \end{pmatrix} \otimes \hat{\rho}_2(t_m) \quad (\text{S10})$$

Starting from a pure state at  $t_0 = 0$ , we have  $(x_0, y_0, z_0) = (x, y, z)$ . The values  $(x_{m+1}, y_{m+1}, z_{m+1})$  depend only on  $(x_m, y_m, z_m)$  and the intermediate idle time  $\tau_m = t_{m+1} - t_m$ : the  $z$  component is perfectly conserved ( $z_{m+1} = z_m$ ), and for the  $x$ - $y$ -components we have

$$x_{m+1} + iy_{m+1} = g(\tau_m) \cdot e^{\pm i\vartheta(\tau_m)} \cdot (x_m + iy_m) \quad (\text{S11})$$

(the “ $\pm$ ” in front of  $\vartheta$  depends on whether the measurement result indicates a rotation into clockwise or counter-clockwise direction) with

$$\tan(\vartheta(\tau)) = e^{-2\tau/T_{\text{triv}}^{(\text{cn})} \cdot \mu_2^2/\mu_1^2} \cdot \sinh\left(4 \cdot \frac{\tau}{T_{\text{triv}}^{(\text{cn})}} \cdot \frac{\mu_2}{\mu_1}\right) \quad (\text{S12a})$$

$$g(\tau) = \frac{1}{\cos \vartheta(\tau)} e^{-2\tau/T_{\text{triv}}^{(\text{cn})}} \quad (\text{S12b})$$

$g(\tau)$  describes the loss in coherence, and  $\vartheta(\tau)$  is the Bayesian guess for the acquired angle. For  $\mathcal{R}_Q$ , we can conclude that

$$\mathcal{R}_Q(t_{m+1}) = g(t_{m+1} - t_m) \mathcal{R}_Q(t_m)$$

In order to determine the optimum value for the idle time  $\tau$ , we compute the effective decay time which is given by  $T_{\text{eff}}(\tau) = -2\tau/\ln(g(\tau))$ ; this form can be obtained easiest by observing the long-term decay for repeated measurements with the same idle time  $\tau$ . Inserting Eq. (S12b), we get

$$T_{\text{eff}}(\tau) = \frac{T_{\text{triv}}^{(\text{cn})}}{1 - \frac{1}{4} \cdot T_{\text{triv}}^{(\text{cn})}/\tau \cdot \ln(1 + \tan^2(\vartheta(\tau)))} \quad (\text{S13})$$

where the explicit form for  $\tan \vartheta(\tau)$  is given in Eq. (S12a).

*1+2 qubit scenario* For one data qubit plus two ancillas, the authors are not aware how to perform an analytical study like for the two-qubit scenario with reasonable effort. Instead, we follow a different strategy: we choose a maximum number of successive measurements within one cycle (we will call this the search depth) to select a subset of all possible feedback schemes, and perform a brute-force search over them.

For a proper definition of the search depth, we start from the realization that if both ancillas are measured without time delay, the quantum system is in a product state of the data qubit and the two ancillas. This means that all correlations are lifted in that process, and thus no further decision can benefit from adaptive response to measurement results before this point in time. Hence, (quasi-)simultaneous measurements split long action sequences into smaller subcycles. For any decision tree (which determines the response to the probabilistic measurement results), we define its depth as the maximum number of idle periods before the cycle is terminated, or equivalently the number of measurements excluding the very last one (that is launched without time delay).

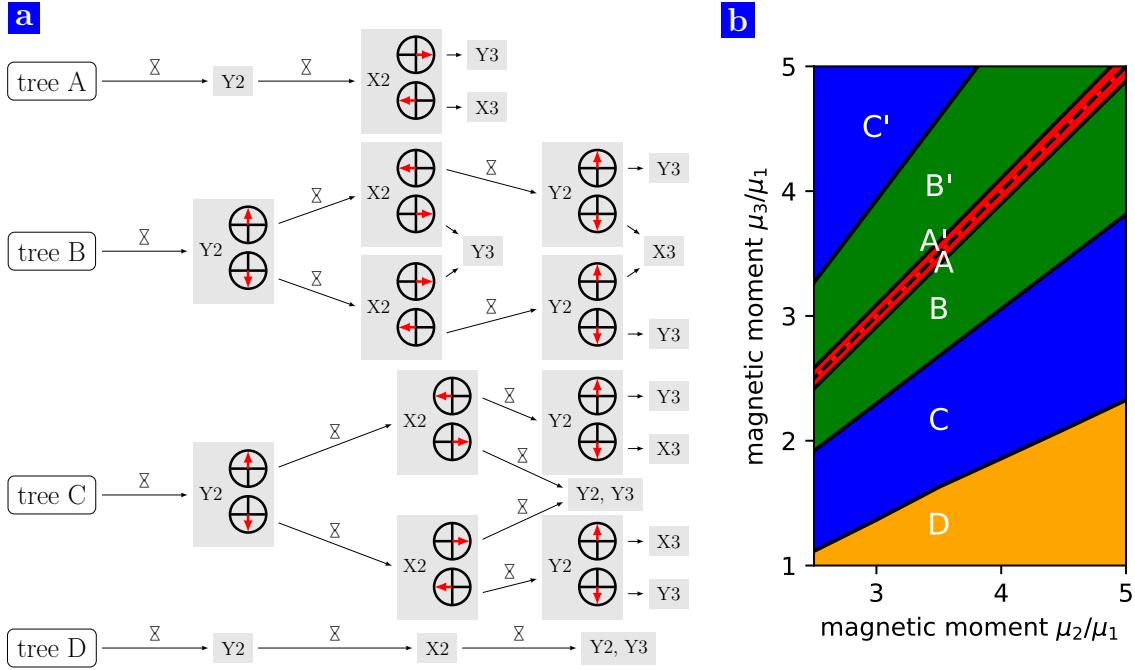


Figure S10. Strategies for the 1+2 qubit correlated noise scenario. (a) Definition of the strategies via decision trees. Idle times are marked with an hourglass symbol. We assume that we always start from a  $+x$  polarization in both qubits. (b) Phase diagram indicating the prevailing strategy for different values for the magnetic moment. X' means that the roles of the qubits 2 and 3 are switched. Note that better strategies might exist which are located outside our restricted search space.

For all strategies that are represented by a finite-depth decision tree, we can compute exact values for effective decay times (in a numerical way, see below). However, there is a considerable limitation for the depth of the brute-force search as discussed in the next paragraph.

The brute-force search does not only have to scan over all (adaptive) decision trees for measurements, but it also needs to find suitable lengths of the intermediate idle times. We discretize the (continuous) parameter range for these idle times by fixing  $n$  different values. For search depth  $d$ , this leads to a total number  $N_d$  of possible feedback schemes following the recursive expression

$$N_1 = 16n \quad (\text{S14a})$$

$$N_{d+1} = 4n \cdot (N_d + 2)^2 \quad (\text{S14b})$$

(derivation see below). From  $N_d \geq 2 \cdot (8n)^{2^d - 1}$ , we can see that there is a double-exponential growth of possibilities with the search depth:

$$N_d \sim \mathcal{O}(\exp(\exp(d))) \quad (\text{S15})$$

$N_d$  grows so fast that in practice only very small search depths  $d \leq 3$  are accessible in reasonable time (but  $d = 3$  only with very high effort). We restrict ourselves to  $d = 2$ , but to still obtain insightful results, we fix the first measurement to the intuitively most reasonable choice (for the ancilla with the largest  $\mu_q$ , the axis orthogonal to its last known position), and due to the symmetry of this situation, we can analyze the two subtrees independent of each other such that the search can effectively be extended by one level. Following this approach, we could identify four different strategies (see Fig. S10a) which prevail for different combinations  $(\mu_1, \mu_2, \mu_3)$  of the magnetic moments (see Fig. S10b). For these four fixed strategies, we in turn run a fine-grid search to further optimize the idle times. Note that all these results might be non-exhaustive due to the various restrictions of the search space.

We proceed by describing the remaining technical aspects of the search technique. First, we will discuss the rationale behind the recursive relation in Eq. (S14):

- $d = 1$ : We can first choose between  $n$  different idle times, and then between 4 different measurement variables (which qubit,  $x$  vs.  $y$ ). For  $d = 1$ , the cycle must be terminated at this point, so an immediate measurement

must follow (the  $x$  or  $y$  axis of the other qubit); this decision may depend on the previous measurement result. In total, this makes  $n \cdot 4 \cdot 2^2$  possible combinations.

- $d \rightarrow d + 1$ : Again, we have the choice between  $n$  different idle times and 4 different measurement variables. Then, we can – dependent on the measurement result – choose between all  $N_d$  combinations in the search tree with depth reduced by 1, or the two instantaneous measurements on the other qubit to terminate the cycle immediately. In total, this makes  $n \cdot 4 \cdot (N_d + 2)^2$  possible combinations for  $N_{d+1}$ .

In order to judge a particular feedback scheme, we have to determine an effective decay time. For each branch  $j$  in the decision tree, we can find out the probability  $p_j$  to end up in this branch, the total time span  $\mathcal{T}_j$  (sum of all idle times), and the loss  $G_j$  of quantum information ( $\mathcal{R}_Q(t + \mathcal{T}_j) = G_j \mathcal{R}_Q(t)$ ). From this information, we can compute an (averaged) effective decay time:

$$T_{\text{eff}} = -\frac{2 \sum_j p_j \mathcal{T}_j}{\sum_j p_j \ln G_j} = -\frac{2 \langle \mathcal{T} \rangle}{\langle \ln G \rangle} \quad (\text{S16})$$

This dependency on  $\mathcal{T}_j$  and  $G_j$  makes it hard to narrow down the search. Suppose  $A = \{j_1^{(A)}, j_2^{(A)}, \dots\}$  and  $B = \{j_1^{(B)}, j_2^{(B)}, \dots\}$  represent two alternatives for the same subtree. If  $\langle \mathcal{T} \rangle_A > \langle \mathcal{T} \rangle_B$  and  $\langle \ln G \rangle_A > \langle \ln G \rangle_B$ , then  $A$  is clearly the better option (one “ $>$ ” and one “ $\geq$ ” would be enough). However, if  $\langle \mathcal{T} \rangle_A > \langle \mathcal{T} \rangle_B$  and  $\langle \ln G \rangle_A < \langle \ln G \rangle_B$ , then it depends on the other branches of the decision tree whether  $A$  or  $B$  should be preferred.

## 7.10. Hyperparameter analysis

In this section, we analyse the influence of hyperparameters in our learning scheme. First, we fix one common set of hyperparameters and apply it to the physical scenarios shown in the main text. In a second step, we modify this hyperparameter set in various ways and discuss the effect of these changes on the learning behavior.

The central results are:

- The learning rate is the only crucial hyperparameter, but there is a clear signature which tells whether it has to be increased or decreased.
- All the other hyperparameters do not decide about the success of the learning process (as long as they stay in reasonable bounds).
- However, some of these hyperparameters have an influence on how fast the learning progresses, and thus can be tuned to optimize the training time.

### a. Common hyperparameter set

We fix one common set of hyperparameters (values are listed in table S3) and apply it to the physical scenarios shown in the main text. The resulting learning curves are plotted in Fig. S11 a. The training works straightforward for all cases except for the “triangle” setup; there, several attempts are required, and we show the best one in each case. However, the limited success rate is not a property of the specific hyperparameter set used here, but also occurs for the hyperparameters used in fig. 5a.

To demonstrate robustness against variations of the physical parameters, we change two of them, the decoherence time (Fig. S11 b) and the number of time steps (Fig. S11 c), but still apply the same hyperparameter set as above. Again, we run training jobs for all the scenarios. Note that in the example with 50 time steps (Fig. S11 c) to keep the number of data points per learning update constant, we have to increase the batch size. Again, learning works smooth for all scenarios, except for the “triangle” setup where multiple attempts are needed.

### b. Modified hyperparameter sets

To analyse how the hyperparameters influence the training progress, we consider several modified hyperparameter sets (table S4) and again apply them to all scenarios (with the original values for decoherence time and number of time steps). Because some of the hyperparameters are implicitly coupled, we have to properly compensate for this. For example, just decreasing the batch size effectively raises the learning rate, and so one would observe the combined

<b>scenarios</b>	
all-to-all	fig. 3 and fig. 4
chain (1 msmt)	fig. 5a (top)
chain (all msmts)	fig. 5a (center)
triangle	fig. 5a (bottom)
all-to-all with small msmt noise	fig. 5c (msmt error probability $\epsilon = 0.01$ )
all-to-all with large msmt noise	fig. 5c (msmt error probability $\epsilon = 0.1$ )

<b>physical parameters</b>	
parameter	a, d) main text values b, e) faster decay c, f) less time steps
a-c) decoherence time $T_{\text{dec}}$ for bit-flip noise	1200 500 1200
d-f) trivial decay time $T_{\text{triv}}^{(\text{cn})}$ for correlated noise	500 200 500
d-f) magnetic moment ratios	2 qubits: $\mu_1 : \mu_2 = 1 : 4$ 3 qubits: $\mu_1 : \mu_2 : \mu_3 = 1 : 3.7 : 4$ 4 qubits: $\mu_1 : \mu_2 : \mu_3 : \mu_4 = 1 : 3.7 : 4 : 4.2$
number of time steps $T$	200 200 50

<b>hyperparameters</b>	
ANN architecture	(# input neurons, 300, 300, # actions)
ANN activation function	Softmax in output layer, elsewhere ReLU
batch size	in a, b, d, e: 64; in c, f: 256
learning algorithm	Adam ( $\eta = 0.0003$ , $\beta_1 = 0.9$ , $\beta_2 = 0.999$ , no bias correction)
# PCA comps in input	6 (in d-f for two-qubit scenarios: 4 because $\dim \mathcal{H} = 4$ )
return discount rate	$\gamma = 0.95$
baseline discount rate	$\kappa = 0.9$
reward scale	$\lambda_{\text{pol}} = 4.0$
punishment reward coefficient $P$	0.1

Table S3. Summary of the physical scenarios, their parameters and the learning hyperparameters as used in Fig. S11. Note that the exact value for the learning rate used in the simulation is in fact  $0.0001\sqrt{10}$ ; the irrational factor of  $\sqrt{10}$  is caused by a slight deviation between our implementation and the standard Adam scheme which in the end resulted only in a redefinition of the learning rate.

name	overridden hyperparameters
low learning rate	Adam learning rate $\eta = 0.00003$
high learning rate	Adam learning rate $\eta = 0.003$
small batch size	batch size 16 Adam hyperparameters $\eta = 0.00008$ , $\beta_1 = 0.975$ , $\beta_2 = 0.99975$
additional layer	ANN architecture (# input neurons, 300, 300, 300, # actions)
pyramid	ANN architecture (# input neurons, 500, 300, 100, # actions)

Table S4. Modified hyperparameter sets as used in Fig. S12. The values given in table S3 are referred to as “standard”, and the right-hand side specifies which values are changed w.r.t. “standard” in the corresponding variation. Note that the Adam hyperparameters are implicitly correlated with the batch size; in “small batch size”, this makes it necessary to adjust  $\eta$ ,  $\beta_1$ ,  $\beta_2$  for compensation to see the direct influence of the batch size (cmp. appendix 7.10.2).

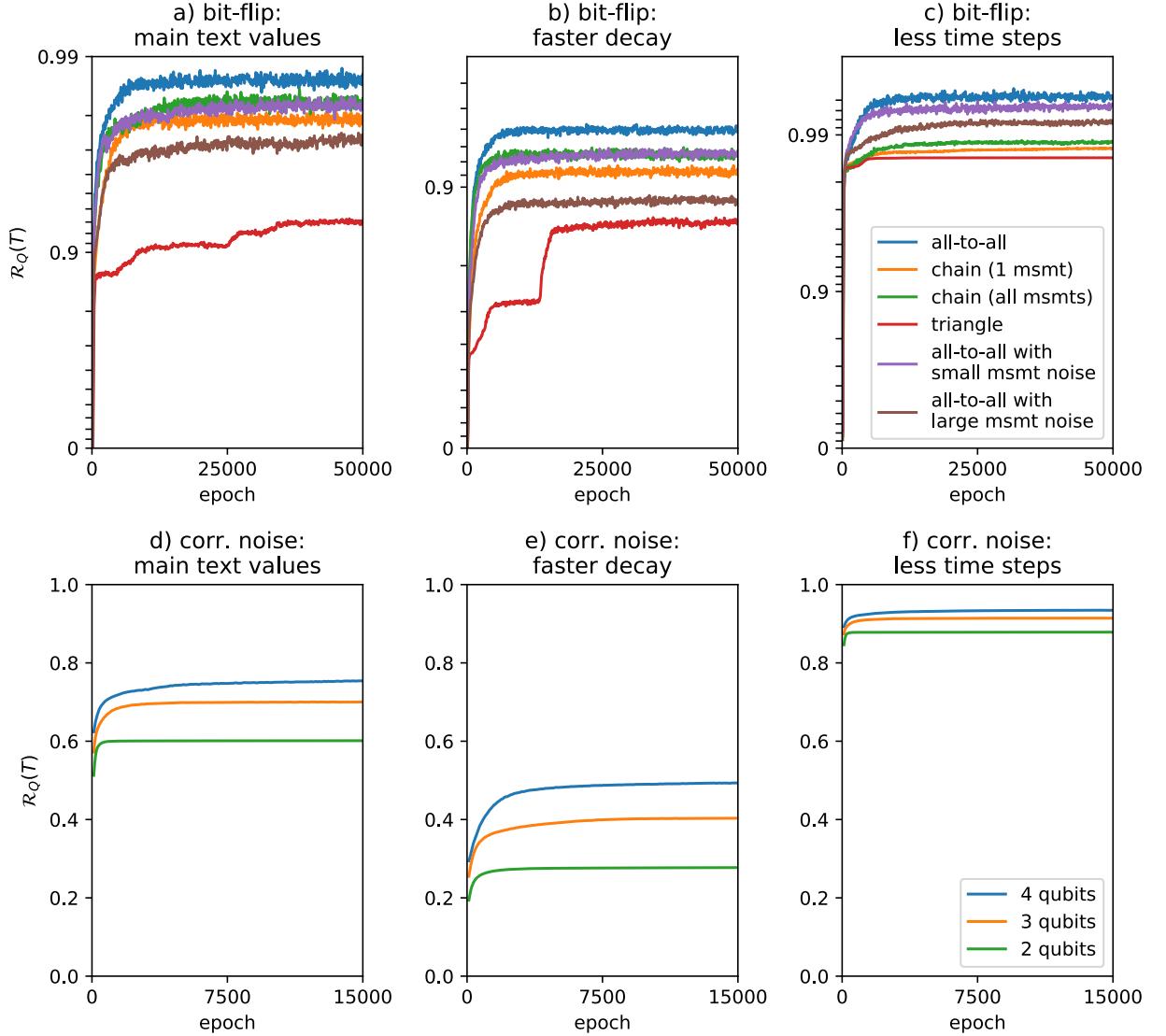


Figure S11. Performance of a common hyperparameter set for different physical systems, defined in table S3. The learning curves show the value of the recoverable quantum information  $R_Q(T)$  at the final time step of the episodes, averaged over the last 100 epochs. We observe successful learning for all the combinations (for the “triangle” scenario, there is a limited success rate, so we trained 5 networks each and show the best run here). Note that the performance for the optimum strategy depends both on the scenario and the physical parameters, so the direct comparison of the saturation levels in these plots against each other does not make a statement about a quality difference in the learning process.

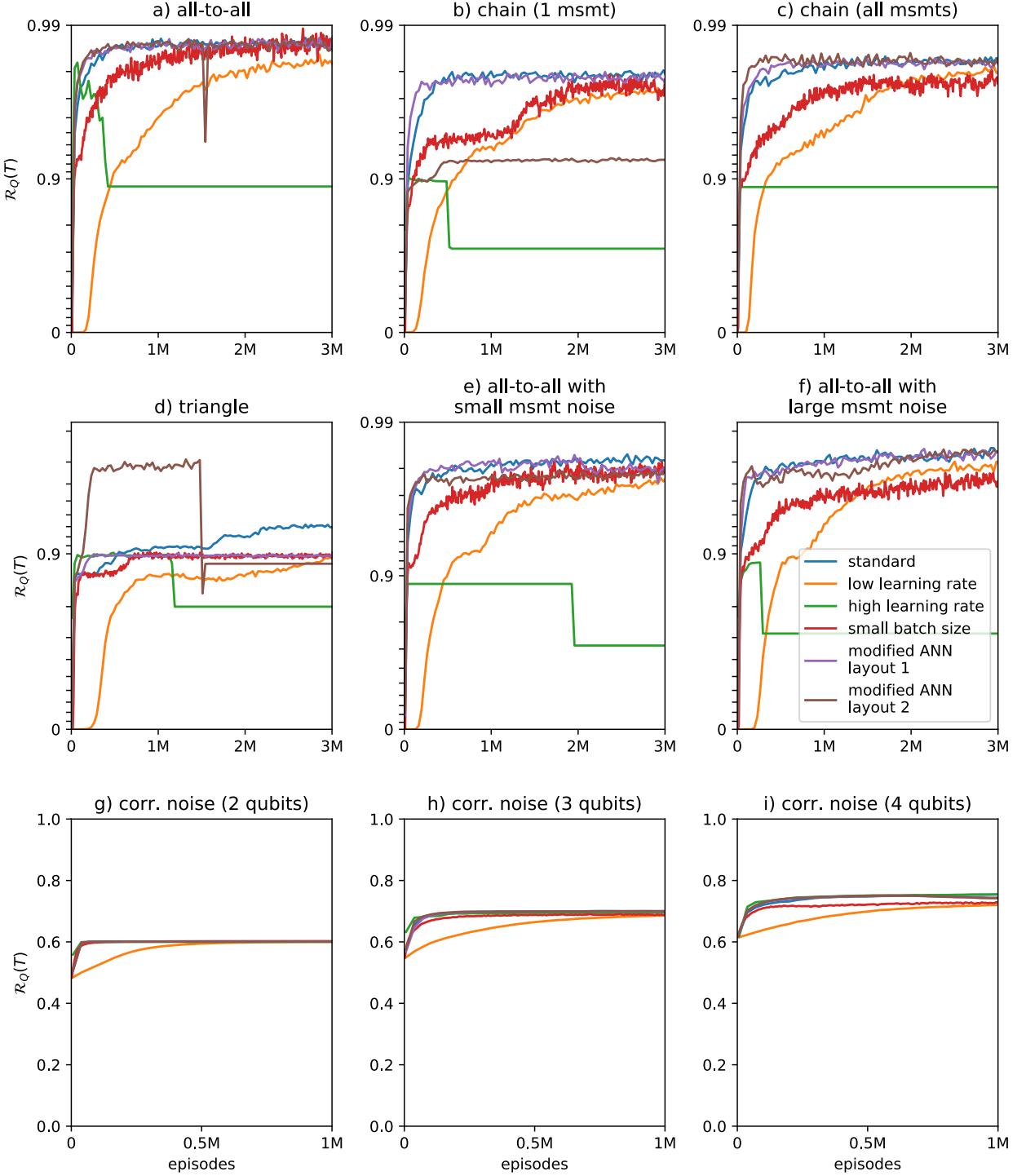


Figure S12. Influence of the hyperparameters on the training behavior. The physical properties are taken from table S3 (main text values), and the different hyperparameter sets are given in table S4. The learning curves show the value of the recoverable quantum information  $R_Q(T)$  at the final time step of the episodes, averaged over the last 100 epochs. The only hyperparameter set which fails (for the bit-flip examples) consistently is “high learning rate”: it starts to learn (faster than the others), but then instabilities show up and finally the training process collapses. All other hyperparameter sets work reliably and reach similar final performance, but not always converge with the same speed. The only exception is the “triangle” scenario: multiple attempts are necessary (independent of the hyperparameter set), so in each case we trained 5 independent networks and show the best run here. In addition, “modified ANN layout 2” for “chain (1 msmt)” (brown curve in b) failed; however, this is an outlier as the majority of the runs performed much better when retraining under the same conditions. The results are discussed in appendix 7.10.2. Note that the x axis displays episodes (epoch times batch size) to allow for direct comparison between runs with different batch sizes.

effects of a smaller batch size and a higher learning rate. We counter these dependencies by suitably adjusting the correlated variables (e.g. here, decreasing the learning rate accordingly), independent of whether this would make a difference or not.

The results are shown in Fig. S12. We make the following observations:

- The learning rate has the largest impact on the learning progress. If it is chosen too small, the convergence is slowed down significantly, whereas if the high learning rate is too high, the situation is even worse as it leads to instabilities which finally lead to a collapse in the learning process. Hence, there is a simple rule of thumb how to find a suitable value for the learning rate: if instabilities show up, the learning rate has to be decreased; as long as the learning curve is smooth, the learning rate can be increased until the point where instabilities occur.
  - The results indicate that a larger batch size seems to be favorable. The most plausible explanation for this behavior is the fact that we use the natural gradient (see appendix 3.3) which involves the Fisher information matrix. Since the exact value of the Fisher information matrix is not accessible, we need to compute an estimate for it (in each epoch), and the quality of this estimate increases with the number of data points, i.e. the batch size. The resulting statistical noise might be amplified by the fact that it is actually the inverse of the Fisher information matrix which enters the calculation of the learning gradient.
  - The network architecture seems to play a minor role.
- 

- [1] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, Cambridge ; New York, anniversary edition edition, January 2011.
- [2] Raymond Laflamme, Cesar Miquel, Juan Pablo Paz, and Wojciech Hubert Zurek. Perfect quantum error correcting code. *Physical Review Letters*, 77(1):198, 1996.
- [3] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [4] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, May 1992.
- [5] D Kingma and Jimmy Ba Adam. A method for stochastic optimisation. In *International Conference for Learning Representations*, volume 6, 2015.
- [6] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- [7] Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276, February 1998.
- [8] Sham M Kakade. A natural policy gradient. In *Advances in neural information processing systems*, pages 1531–1538, 2002.
- [9] Jan Peters, Sethu Vijayakumar, and Stefan Schaal. Reinforcement learning for humanoid robotics. In *Proceedings of the third IEEE-RAS international conference on humanoid robots*, pages 1–20, 2003.
- [10] Razvan Pascanu and Yoshua Bengio. Natural gradient revisited. 2013.
- [11] Ronald J Williams and Jing Peng. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.
- [12] Sepp Hochreiter and Jürgen Schmidhuber. Long-Short Term Memory. *Neural Computation*, 9:1735, 1997.
- [13] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [14] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579, 2008.
- [15] Rami Al-Rfou et al. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688, May 2016.
- [16] Daniel Gottesman. The heisenberg representation of quantum computers. *arXiv preprint quant-ph/9807006*, 1998.
- [17] Richard Cleve and Daniel Gottesman. Efficient computations of encodings for quantum error correction. *Physical Review A*, 56(1):76, 1997.
- [18] Scott Aaronson and Daniel Gottesman. Improved simulation of stabilizer circuits. *Physical Review A*, 70(5):052328, 2004.
- [19] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825, 2011.