

Quantum Science and Technology

OPEN ACCESS

CrossMark

PAPER

Randomized benchmarking in the analogue setting

RECEIVED

1 October 2019

REVISED

25 February 2020

ACCEPTED FOR PUBLICATION

11 March 2020

PUBLISHED

23 April 2020

E Derbyshire^{1,4} , J Yago Malo², A J Daley², E Kashefi^{1,3} and P Wallden¹ ¹ School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, United Kingdom² Department of Physics and SUPA, University of Strathclyde, Glasgow G4 0NG, Scotland, United Kingdom³ Laboratoire d'Informatique de Paris 6, CNRS, Sorbonne Université, 4 place Jussieu, 75005 Paris, France⁴ Author to whom any correspondence should be addressed.E-mail: E.E.Derbyshire@sms.ed.ac.uk**Keywords:** analogue quantum simulation, randomized benchmarking, approximate unitary t -designs

Original content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](#).

Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.



Abstract

Current development in *programmable* analogue quantum simulators (AQS), whose physical implementation can be realised in the near-term compared to those of large-scale digital quantum computers, highlights the need for robust testing techniques in analogue platforms. Methods to properly certify or benchmark AQS should be efficiently scalable, and also provide a way to deal with errors from state preparation and measurement (SPAM). Up to now, attempts to address this combination of requirements have generally relied on model-specific properties. We put forward a new approach, applying a well-known digital noise characterisation technique called randomized benchmarking (RB) to the analogue setting. RB is a scalable experimental technique that provides a measure of the average error-rate of a gate-set on a quantum hardware, incorporating SPAM errors. We present the original form of digital RB, the necessary alterations to translate it to the analogue setting and introduce the analogue randomized benchmarking protocol (ARB). In ARB we measure the average error-rate per time evolution of a family of Hamiltonians and we illustrate this protocol with two case-studies of analogue models; classically simulating the system by incorporating several physically motivated noise scenarios. We find that for the noise models tested, the data fit with the theoretical predictions and we gain values for the average error rate for differing unitary sets. We compare our protocol with other relevant RB methods, where both advantages (physically motivated unitaries) and disadvantages (difficulty in reversing the time-evolution) are discussed.

1. Introduction

In the quest for real-world applications in the present and near-future, the focus of quantum computing has fallen on noisy intermediate-scale quantum systems. While there is a lot of interest in digital quantum computing/simulation (DQC), significant progress has been made in analogue quantum simulation (AQS), which offers existing medium-scale systems, and where relevant physical problems are mimicked by a highly-tunable quantum system. AQS are engineered to evolve continuously in time according to a specific Hamiltonian or family of Hamiltonians which can be flexibly controlled and microscopically understood from first principles. Instead of the application of quantum logic gates, as in the digital setting, a calculation is performed through the unitary or dissipative evolution under the system Hamiltonian. As a result, AQS currently lacks many of the existing testing or error-correcting methods of DQC. Being able to trust that AQS are simulating the designated quantum system, or running the same class of Hamiltonians in a *reproducible* way becomes essential when the AQS cannot be classically simulated.

Most AQS experiments aim to estimate the ground state of a Hamiltonian or to learn about the dynamical properties of the system and the spread of information, i.e. on a given quantum hardware or when evolving under a given Hamiltonian. Typically, AQS is tested against accurate numerical methods in classically solvable regimes (such as one-dimensional dynamics) or cases where the final outcome is known, i.e. classical optimization. More recently, a number of ideas have been developed for tackling classically

intractable regimes. These include a self-verifying technique for simulations of the lattice-Schwinger model [1], a method based on non-demolition measurements of the Hamiltonian [2] and the broader concept of cross-platform verification, i.e. comparing results for the same problem from different quantum hardware or the measurement of compatible correlation functions that can verify specific Hamiltonians even when the quantum dynamics *cannot* be classically simulated [3, 4]. It is also possible to estimate the final state of the AQS through quantum process/state tomography (QPT/QST) [5, 6], direct fidelity estimation (DFE) [7], or in recent years, matrix product state (MPS) tomography [8] and neural-network approaches [9].

QST is a technique used to reconstruct an unknown quantum state when given multiple copies of it to measure, and to estimate the final state of a quantum simulation, whilst QPT is a similar technique used to reconstruct an unknown quantum process (quantum channel) and to test the process of a quantum simulation. Both techniques require resources that grow exponentially with the size of the system. DFE is a technique to determine whether a system arrives at some target state or runs a target gate. This is more efficient than the previous techniques since knowing the target state/operation means that there are fewer resources (measurement settings) required, however it does not directly account for the errors from state preparation and measurement (SPAM) and is still not scalable in the sense that the *actual* number of measurements still scales exponentially. MPS tomography provides an estimate of the final state of the system, and also requires less measurement settings than standard tomography due to the fact that it exploits tensor network techniques (and specifically matrix product states) to approximate the final state; again, it does not include SPAM errors in its analysis and in the suitable non-classical regime, scales exponentially in terms of resources. To summarise, (i) none of the techniques mentioned previously succeed in giving an account of the errors from state preparation and measurement, meaning that the accuracy to which they characterise the noise of the system is always bounded by an unknown contribution from the SPAM errors and (ii) these methods are not efficiently scalable, and can only provide characterisation for systems of up to ~ 20 qubits failing to capture correlations further than a few sites [8] or relying on state symmetries [10].

Despite the considerable advances in recent years [1] most of the current methods rely on model-specific properties to improve on previous proposals. The aim of our research is not to verify the ground state estimate of complex Hamiltonians, but instead find a general way to capture the performance of an analogue quantum simulator in running a family of Hamiltonians. Randomized benchmarking (RB) is a digital method to find a measure for the holistic performance of a quantum hardware that takes into account SPAM errors and is theoretically efficiently scalable, i.e. in the length of the gate sequences applied ($\text{poly}(L)$) and in the size of the system ($\text{poly}(N)$). This scalability is dependent on the efficiency with which the gate-set tested can be composed/compiled with the native operations of the chosen hardware.

The motivation for this work grew from the need to overcome the limitations of previous methodologies and develop scalable non-model-dependent methods, building on ideas from the digital setting and applying them to AQS, i.e. with the intention to reduce the gap in testing techniques that exists between the digital and analogue fields of quantum computing. RB in the analogue setting would offer a way to determine whether your chosen quantum architecture will reliably perform a set of quantum evolutions or unitaries, giving an average error rate for this set of unitaries; particularly useful when considering *programmable* analogue quantum simulation. RB has the potential to efficiently provide a measure of performance of AQS in regimes that are currently classically intractable without the SPAM error noise floor on the accuracy of this partial noise characterisation. Moreover, it will become apparent that RB is more natural and physically motivated for analogue systems than for their digital counterparts. It is also important to mention that the method presented below particularly in the case of perfect time inversion is similar to the Loschmidt echo [11, 12], a well-known technique relevant in the context of quantum chaos.

With these motivations in mind, we propose extending randomized benchmarking to the analogue setting, which we call *analogue randomized benchmarking* (ARB). In section 2 we present the original form of randomized benchmarking and the technical details of it, in order to keep this work self-contained and to better explain the alterations made for our contribution. Following in section 3, the modifications necessary to extend randomized benchmarking to analogue quantum simulators; including the current form of the analogue randomized benchmarking (ARB) protocol (section 3.3). In section 4 we first present two case studies of simulating ARB on concrete analogue models (section 4.1) and then in section 4.2 we analyse the robustness of this characterisation in the presence of more complex noise, dissipation and experimental imperfections. Note that we do not perform experiments on physical hardware, but instead emulate these devices by classically modelling physically motivated noise scenarios. We conclude in section 5 with considering the barriers for physical implementation of this protocol as well as some ideas to overcome them, and introducing the possible future directions that have transpired through this process.

2. Randomized benchmarking

Randomized benchmarking (RB) [13–16] is a technique for evaluating the performance of a quantum hardware which simplifies the error channel of a quantum process such that one can extract the average error per gate of a particular gate-set on this hardware. This measure is for each gate when it is run as part of a long random computation, a relevant metric for most quantum algorithms that rely on this type of computation. In its standard form, it relies on the gate-set having a particular distribution, and on being able to efficiently invert a string of gates from the gate-set with one single operator. A general randomized benchmarking protocol will consist of variants of the following steps: (1) preparing an initial state $\rho_\psi = |\psi\rangle\langle\psi|$ on a quantum device, (2) running very many random sequences of gates of varying lengths such that the system should return to this initial state, (3) measuring the probability that the state remained unchanged, and lastly, (4) plotting and fitting the results to a pre-determined decay model that characterises the average error rate of the gates.

The central theory is based on two observations, the first is that *any* error channel, irrespective of the correlations it initially exhibits, when *twirled* (essentially averaged over random unitaries, see appendix A) ‘behaves’ as a much simpler error channel that is easier to characterise and quantify. Twirling is the process that transforms a quantum channel $\Lambda(\rho)$ into twirled channel $\Lambda_t(\rho)$ by conjugating over unitaries $U(\rho) \in U(d)$, in the following way:

$$\Lambda_t(\rho) = \int_U d\mu(U) U \circ \Lambda \circ U^\dagger(\rho), \quad (1)$$

where d is the dimension of the Hilbert space. If the unitaries $U(\rho)$ are distributed according to the Haar measure $d\mu$ [17], a measure of uniformity, then the twirled channel becomes a *depolarising channel*. A depolarising channel is a simple quantum channel of the following form: $\mathcal{E}_d(\rho) = p_\mathcal{E}\rho + (1 - p_\mathcal{E})\frac{I}{d}$. Intuitively, with some probability $p_\mathcal{E}$ the depolarising channel $\mathcal{E}_d(\rho)$ leaves the state ρ intact, whilst with the remaining probability $(1 - p_\mathcal{E})$ returns the maximally mixed state $\frac{I}{d}$; where I is the identity operator. One of the first papers to suggest this kind of method, by Emerson *et al* [13], used the equivalence of a Haar-twirled quantum channel to its average fidelity over Haar-random unitaries, to construct a motion reversal protocol which results in a single parameter (depolarising parameter) to describe a noisy quantum channel. This was introduced as a way to efficiently estimate the strength of the noise channel on a device and with this method, errors from state preparation and measurement (SPAM), which are independent of the length of a sequence, can also be extracted from the error characterisation. However, this early model had the obstacle that implementing sequences of Haar-random unitaries is inefficient.

This brings us onto the second key observation, which is that the averaging over this infinite (Haar) set of unitaries can be mimicked by sampling from a small finite set. Finite sets that approximate this average are known as *unitary t -designs*; randomly sampling from a unitary t -design is equivalent to randomly sampling from the Haar random unitaries provided that the *averaged* quantity computed involves polynomials of order, at most, t . Formally, a unitary t -design is a set of unitaries $\{U_k\}$, where $\{k = 1, \dots, K\}$ such that:

$$\frac{1}{K} \sum_{k=1}^K P_{t,t}(U_k) = \int_{U(d)} d\mu(U) P_{t,t}(U), \quad (2)$$

for every polynomial $P_{t,t}(U)$ of order t , where $d\mu(U)$ is the Haar distribution. In other words, for any polynomial of unitaries of degree t or less, calculating the average over the set $\{U_k\}$ is the same as calculating the average over all the unitaries (Haar integral). An ϵ -approximate t -design is a set of unitaries that has the same property, where it holds only up to some error ϵ . For RB, the gate-set must be *at least* an exact or approximate 2-design (for intuition as to why this is the case, please see appendix C). In summary, RB aims to test and quantify how well a set $\{U_k\}$ performs on average on a given quantum hardware, by utilising the twirling property of the Haar distribution. The twirling property is then used to simplify the effective error channel and make it quantifiable with a single parameter (depolarisation). In other words, if one has a set that is a unitary t -design, then by performing the method of randomized benchmarking one obtains an average error rate of the said (gate)-set on that quantum device.

We present here a basic randomized benchmarking protocol (see algorithm 1), for completeness, and for technical details and further explanation of how the method works we refer to appendix A. In this protocol, Λ_{U_k} is the imperfect implementation of U_k , ρ_ψ is the initial state, taking into account preparation errors, whilst E_ψ is the positive-operator valued measure (POVM) element taking into account measurement errors; in the ideal case $E_\psi = \rho_\psi = |\psi\rangle\langle\psi|$. The probability of the initial state surviving the random sequences is called the *survival probability* where for a random quantum circuit U_C , $P := \langle\psi|\Lambda_{U_C}(\rho_\psi)|\psi\rangle = \text{Tr}(E_\psi\Lambda_{U_C}(\rho_\psi))$, with Λ_{U_C} the imperfect implementation of U_C . The *average survival*

Algorithm 1. Analogue randomized benchmarking.

-
- 1: Sample uniformly from $\{U_k\}$ a number of sequence time-lengths S_T and run a sequence S_η of time-length $T \in S_T$ where: $S_\eta = [\Lambda_{U_{k_1}}, \dots, \Lambda_{U_{k_\eta}}, \Lambda_{U_{k_\eta}}^\dagger, \dots, \Lambda_{U_{k_1}}^\dagger]$ on your system such that if each unitary was perfectly implemented your system would be returned to initial state ρ_ψ . Here we systematically invert each preceding unitary.
 - 2: Repeat the sequence R times; record the survival probability $P_\eta = \text{Tr}[E_\psi S_\eta(\rho_\psi)]$ for this sequence.
 - 3: Repeat the above for various sequences of the same time-length to get the average survival probability for each sequence time-length: $P_T = \text{Tr}[E_\psi S_T(\rho_\psi)]$, where S_T is the average over all sequences of time-length, T .
 - 4: Repeat the above steps for sequences of different time-lengths and plot the average survival probability against time-length P_T vs T .
 - 5: Fit the results to a predetermined decay curve of the following or similar form: $P_T = A + Bf^T$, where again T is the sequence time-length and f represents the fidelity decay parameter of the process, with A and B fit parameters that absorb SPAM errors. And again, the average error rate is characterised by r where $r = (d-1)(1-f)/d$ and d is the dimension of the Hilbert space for a system of qubit size n ($d = 2^n$).
-

probability (P_l) in the protocol, is equivalent to a product of twirled depolarising channels and due to the left-invariance of the Haar-twirl, can be compared to the average fidelity of the error channel (see appendix A) which results in the decay curve: $P_l = A + Bf^l$. Therefore when the data is fit to the curve, this gives an average error rate r , where $r = 1 - F_{\text{ave}} \equiv (d-1)(1-f)/d$, where F_{ave} is the average fidelity of the gates in the gate-set. The fit parameters A and B vary for different versions of the protocol, and in the simplest case are $A = \frac{1}{d}$, $B = \frac{d-1}{d}$.

In order for the above protocol to really quantify the average error-rate of a gate-set, as given in step 5, we need to make a number of simplifying assumptions on the type of noise (error-channels) the device has: the errors should be (i) gate-independent $\Lambda_{U_k} = \Lambda$ and (ii) time-independent, i.e. independent of the time it takes to run the gate and of when it is applied in any part of any sequence, (iii) the error-channel should be trace-preserving and memoryless (iv) the SPAM errors should be length-independent and (v) the error in the inversion step can be viewed as a single step; therefore, it does not scale with the length of the sequence which allows it to be absorbed into the SPAM errors. Finally, given that the gate-set tested is not universal, it is also crucial that (vi) the inversion step (in the ideal case) can also be implemented using gates from the tested gate-set. Realistically, the physical errors on a quantum process are likely to be gate-dependent and time-dependent, and more recent protocols show RB to be robust against certain types of noise [18–24], however the method is most straightforward when the noise assumptions are adhered to. In our investigation, therefore, we first analyse the simplest case (of time and gate-independent noise, section 4.1) and build up our protocol with different noise models in section 4.2.

The Clifford group of operators is a well-known and simple to construct 2-design [25], and the average performance of these gates is a relevant parameter for several reasons, including error-correcting codes, and the fact that with the addition of one extra single-qubit unitary gate, the set becomes universal. Therefore, many RB protocols test the Clifford group, utilising the fact that a single inversion operator may be found efficiently for the Clifford group [26] such that the protocol consists of running strings of random quantum gates where the last gate is the inversion operator. Efficiently finding a single inversion operator in the analogue setting is currently not possible. Due to this, we adopt a method similar to Emerson *et al*'s [13] debut of running imperfect unitaries and their inverses for our analogue randomized benchmarking protocol, rather than a single inverse for a string of gates, where ours is closer to the form of the Loschmidt echo sequence [11, 12]. Our protocol differs in that we do not test Haar-random unitaries nor even an exact 2-design, instead aiming to benchmark an ϵ -approximate 2-design. There are several RB methods that are relevant to our exploration, which we highlight in the next section.

3. The analogue setting

In this section, we first give an overview of extending RB to the analogue setting, with technical details from section 3.1 onwards, and our ARB protocol (see algorithm 2) in section 3.3. We replace the quantum logic gates tested by digital RB (most commonly, the Clifford group or generators thereof) with a set of time evolution operators/unitaries: $\{U_k = e^{-iH_k dt}\}$. One of the key conditions of performing RB with these types of unitaries is that they must converge to an approximate unitary 2-design, i.e. form an ϵ -approximate 2-design. There have been interesting results in benchmarking finite groups that are not the Clifford group or a known 2-design [22, 23, 27], in particular direct [23] and approximate [22] randomized benchmarking; in the former, rather than benchmarking the full Clifford group the protocol tests a set of ‘native gates’ with the requirement that they generate the Clifford group. The latter writes the RB protocol in terms of arbitrary finite groups giving a bound on results from RB when sampling from an approximate Haar-distribution on this group compared to sampling from the full Haar-distribution. Our protocol is similar in the sense that we aim to sample from an approximate Haar distribution, using a finite set of

Algorithm 2. Digital randomized benchmarking.

-
- 1: Sample uniformly from $\{U_k\}$ a number of sequence lengths S_l and run a sequence S_η at length $l \in S_l$ where:
 $S_\eta = \Lambda_{U_{k\eta+1}} \Lambda_{U_{k_1}}, \dots, \Lambda_{U_{k\eta}}$, and $\Lambda_{U_{k\eta+1}}$ is a single operator deterministically chosen to invert the preceding sequence of unitaries (i.e. $\Lambda_{U_{k\eta+1}} = [\Lambda_{U_{k_1}}, \dots, \Lambda_{U_{k\eta}}]^\dagger$). This sequence should return the system to its initial state ρ_ψ .
 - 2: Repeat this sequence R times and record $\text{Tr}[E_\psi S_\eta(\rho_\psi)]$ to see if initial state ρ_ψ survived the sequence S_η and call this the survival probability P_η for sequence η .
 - 3: Repeat this for varying sequences of the same length l and find the average probability that the initial state survived for this sequence length, $\text{Tr}[E_\psi S_l(\rho_\psi)]$ where S_l represents the average over all sequences of length l . Call this the average survival probability for length l : P_l .
 - 4: Repeat the above steps for sequences of different lengths, and plot average survival probability against sequence length, i.e. P_l vs l . Fit results to a pre-determined decay curve: $P_l = A + Bf^l$ where l is the sequence length and f is the fidelity decay parameter, with A and B absorbing SPAM errors.
 - 5: The average error rate can be characterised by r where $r = (d-1)(1-f)/d$, and d is the dimension of the Hilbert space for a system of qubit size n ($d = 2^n$).
-

unitaries, however we do not require that this set forms a group, we only require that this set approximates a 2-design. Therefore, we gain some freedom in our approach which allows us to test more quantum hardware that can not fulfill this stronger condition, such as analogue quantum simulators. The result of [22] is also important to our protocol, because the authors demonstrate that equally meaningful results can be gained from RB with a distribution close enough to that of a 2-design (or fully Haar-random).

We create a set of unitaries from a disordered set of multi-qubit Hamiltonians (details to follow), in a way that is more easily implementable on a physical device, and we sample unitaries from this set to create long sequences such that we assume after some time that we approximate a 2-design. We view our unitaries as generators of an approximate 2-design and our aim is firstly, to demonstrate heuristically that the disorder we create in our Hamiltonians does in fact result in convergence to an approximate 2-design and secondly, that this convergence rate is sufficient enough for approximate twirling (through ARB) to characterise realistic noise in this setting. Furthermore, no RB scheme has been attempted in the analogue setting nor directly adapted to the natural evolution of a quantum device. One of the caveats of digital one and two-qubit gate RB methods is that they require compilation of physical gates, that are produced naturally in the device, into gates of the Clifford group or other digital gates that generate such a group. Compilation of *any* logical gates that are not native to the device limits the size of the system that one can feasibly test, and we bypass this compilation issue by using only the native capabilities of the device.

3.1. Unitary set for ARB

The unitaries tested are created by disordering a Hamiltonian H_s , native to the quantum device of interest, resulting in a set of Hamiltonians $\{H_k\}$ that when time-evolved produce a set of unitaries $\{U_k = e^{-iH_k dt}\}$. Our aim is to construct a technique that while general, takes into account the physical limitations and the possibilities offered by existing systems. For this reason, the choice of dt is determined by factoring in physical constraints of the device; it should not be less than the minimum time required to *physically allow* for the Hamiltonian to be changed per time-step. We define H_k to be:

$$H_k = H_s + \zeta_k^{(g,l)}, \quad (3)$$

where $\zeta_k^{(g,l)}$ is an added disorder term which we define to be one of the following:

$$\begin{aligned} \zeta_k^g &= \Delta_k \sum_{ij} \sigma_i^u \otimes \sigma_j^u \\ \zeta_k^l &= \sum_{ij} \Delta_k^{ij} \sigma_i^u \otimes \sigma_j^u, \end{aligned} \quad (4)$$

where the indexes (g, l) denote global (g) and local (l) disorder terms, and $\Delta_k^{(ij)}$ is a disorder potential that varies for every $\{k = 1, \dots, K\}$ and, in the case of local disorder, according to which sites it is acting on. Here the index $u = x, y, z$ indicates which product of Pauli operators is to be applied on nearest-neighbours. By varying the original Hamiltonian H_s with these disorder terms we generate a family of Hamiltonians. From this family $\{H_k\}$, we generate the time evolution operators to obtain our unitary set $\{U_k\}$ for a fixed time-step (dt).

Producing an ϵ -approximate 2-design in this setting is an interesting study, and there are relevant results demonstrating convergence to 2-designs for time-dependent Hamiltonians exhibiting Brownian motion [28], for locally disordered Hamiltonians [29] and the use of locally random unitaries to estimate Rényi entropies [30]. In the latter study of Rényi entropies the authors found that by applying a local disorder

potential and a single-site Pauli-operator to all sites of a Hamiltonian sector ($H|_A$) at each time-step, their random unitaries converged to a unitary 2-design on the sectors, defined from the irreducible representation decomposition of the group generated by their unitaries. Our Hamiltonians, on the other hand, are time-independent and the aim is to converge to a 2-design over the entire Hilbert space of our system. Their results did however influence our choice of disorder potential and we apply disorder potentials also drawn from a normal distribution, with standard deviation $\delta = J$; however, we apply these potentials globally (homogeneous along the chain) and locally, in order to generate two different sets of unitaries for each model, applying them with a standard spin–spin interaction term. For a fixed Hamiltonian H_s , the generated time-evolution operator can never approach a 1-design, let alone a 2-design [31], because the distribution of the eigenvalues will be too localised. Disordered Hamiltonians have a rich body of literature exploring their role in quantum thermalization, quantum chaos, scrambling, and random unitaries [32–37]. Intuitively, for a set of Hamiltonians of the form of equation (3) to converge to a unitary 2-design, the disorder term should not conserve any of the symmetries of the Hamiltonian H_s . If the disorder term were to commute with parts of H_s and therefore conserve, for example, the total spin, then the disordered Hamiltonians would be exploring only subspaces of the Hilbert space [38]; whereas, breaking all the symmetries of a Hamiltonian should produce statistics associated with random unitaries [29, 31, 34, 39–41]. Initially, we expected that our sets of unitaries would in themselves be ϵ -approximate 2-designs, and with sufficiently large K we predict this to be the case. For all of our simulations, however, we fix $K = 1000$, as a large finite set of unitaries that is experimentally feasible. Subsequent investigation indicated that the combination of the elements of these sets produce an approximate 2-design after some time (with some sets performing better than others, see section 4) hinting that our sets are, rather, generators of ϵ -approximate 2-designs.

In ARB the main operation that we want to achieve is twirling the error channel into a depolarising channel by averaging this channel over our unitary set (see section 2 and appendix A); with the steps in standard RB performing this necessary twirl because the unitaries that are randomly sampled are *uniformly* distributed. As mentioned, there is evidence that approximating a 2-design is adequate for RB [22] however there is no guarantee that the errors will be depolarised with an ϵ -approximate 2-design since the unitaries will not come from an *exact* uniform distribution. The value of ϵ is crucial, with only a sufficiently good approximate design producing meaningful results from RB. We assume that our sets produce an ϵ -approximate 2-design for long sequences, and we derive bounds on our results in section 4 based on this assumption. Here, we define an ϵ -approximate 2-design and present these bounds.

Definition 3.1. An ϵ -approximate unitary 2-design defined in terms of the diamond-norm [42] is a measure on a finite subset $U(D)$, where D is the dimension of the Hilbert space, that satisfies the following property [43]:

$$\|E_\alpha(\Lambda) - E_\mu(\Lambda)\|_\diamond \leq \epsilon, \quad (5)$$

where E_α is the twirled channel of Λ over a set of unitaries $\{U_\alpha\}$ spread according to a probability distribution α . E_μ is the Haar-twirl of that channel, with μ the Haar measure.

A unitary 2-design corresponds to the case when $E_\alpha = E_\mu$. There are a few ways to determine how close a particular set of unitaries is to an exact unitary 2-design, that involve comparison with the way that Haar random unitaries behave, such as the frame potential [44] and comparing the second moment operators [45] (see appendix F).

Theorem 3.1. Let $P_l^\mu = A + Bf^l$, $P_l^\alpha = P_l^\mu \pm \delta P_l$ be the average survival probabilities measured from randomized benchmarking for a sequence length, l , of unitaries distributed according to the Haar measure, μ , and a set of unitaries, $\{U_\alpha\}$, distributed according to an unknown α , respectively. If the set $\{U_\alpha\}$ is an ϵ -approximate 2-design, it holds that:

$$|P_l^\alpha - P_l^\mu| = |\delta P_l| \leq l \cdot \epsilon, \quad (6)$$

where $P_l^\tau = \text{Tr}[E_\psi E_\tau(\Lambda)^l(\rho_\psi)]$, $\tau \in \{\alpha, \mu\}$ is the survival probability of input state ρ_ψ when the twirled quantum channel E_τ is applied to it l times.

Proof sketch. To go from the *measured* survival probabilities to the exact case, one needs to replace the l sums over the ϵ -approximate 2-design with the corresponding integrals over the Haar measure. For each one of these replacements, the maximum difference between the two probabilities increases by ϵ , resulting in the $l \cdot \epsilon$ at the end of theorem 3.1. The full proof is given in appendix F. \square

Assumption 3.1. We assume that the statistical error in estimating the value of P_l^α with the RB method, is much smaller⁴ than the error induced by running RB with an ϵ -approximate 2-design rather than an exact

⁵ This assumption can always be fulfilled with a suitable choice of the number of repetitions R of each sequence.

2-design:

$$\delta P_l^\alpha \ll \delta P_l. \quad (7)$$

Theorem 3.2. *Let A and B be known quantities, and let r' be the average error rate determined from RB for the set of unitaries $\{U_\alpha\}$ that are an ϵ -approximate 2-design. Where the average error rate for Haar-distributed unitaries $\{U_\mu\}$ is $r = (d-1)(1-f)$, and $d = 2^n$ where n is the dimension of the Hilbert space the unitaries are applied to. Under assumption 3.1 where theorem 3.1 holds, with $l = 1$ and no errors from state preparation or measurement, r' is bounded as follows:*

$$r - \epsilon \leq r' \leq r + \epsilon. \quad (8)$$

The proof of theorem 3.2 can be found in appendix F.

3.2. Systematically inverting unitaries

The standard form of RB (see algorithm 1) involves a single deterministically chosen inversion operator $\Lambda_{U_{k+1}}$ that inverts the preceding unitary sequence. The errors in the process before this inversion step get depolarised through the twirling of the error channel, but the single inversion operator will also have errors attached to it. Due to this inversion operator being much shorter in length than the sequence preceding it, the error associated with this operator is a constant additive error. While this (SPAM & inversion) error is not (in general) depolarising, as far as computing the survival probabilities, there exists some depolarising error channel that would have the exact same effect⁶, and therefore we can model it as such.

As mentioned, the inversion operator of any string of Clifford operations can be found efficiently. Unfortunately, there is no equivalent result in the analogue setting, and therefore the initial development of the protocol involves systematically inverting each preceding unitary, like the Loschmidt echo sequence [11, 12], and similar to [13]. This systematic inversion of the preceding sequence means that the errors, now a combination of forward evolution and inversion errors do not, in general, depolarise in the RB process; therefore, for the purposes of analysis in section 4.1, we model the systematically inverted unitaries as perfect. We explore the scenario of having the same type of noise in the inversion operators in section 4.2. A necessary step to physically implementing analogue randomized benchmarking will be to remove the time-reversal aspect of our protocol. Firstly, because with systematic inversion the errors will not necessarily be depolarised, and secondly and, most importantly, because performing this type of time-reversal is not trivial on the analogue quantum simulator. While certain terms such as field or on-site terms can be inverted efficiently in the current experiments, the inversion of off-site couplings or tunneling terms still remain challenging [46, 47]. To highlight this non-triviality, we note that only recently was such inversion realised even in a simplified model [48].

3.3. Analogue randomized benchmarking protocol

Determine a set of Hamiltonians $\{H_k\}$ with $k = \{1, \dots, K\}$ such that $\{U_k\}$ forms a sufficiently good ϵ -approximate unitary 2-design for twirling over them to approximate a depolarising channel, where $U_k = e^{-iH_k dt}$. The time-step dt is kept the same for each unitary operator to mitigate time-dependent errors (except in section 4.2). We again prepare the quantum system in initial pure state ρ_ψ where, if ideally prepared, $\rho_\psi = |\psi\rangle\langle\psi|$ and assume that the error channel is a trace-preserving and memoryless CPTP map, with errors being gate and time-independent.

The parameters for the standard RB protocol (see algorithm 1 in section 2) still apply, but we redefine the form of S_l and introduce another measure for the sequence length: S_T where $S_T \equiv S_l$ and T represents the total time to run each sequence of length l , i.e. $T = dt \cdot l$ and so, we highlight that the average error-rate gained from this protocol is the average error as a function of physical time T . Now, our sequence lengths are called S_T and we refer to length as time-length.

In the next section we focus on describing how ARB would operate on an analogue device, where our analysis involves classically simulating the process, including modelling and implementing the error channels.

4. Case studies

To analyse ARB, we consider a spin system native to many quantum simulators, which is particularly relevant to the case of trapped ion experiments [49–51]. Specifically, we modelled the XY Hamiltonian:

⁶ This happens because we measure in the basis $\{|\psi\rangle\langle\psi|, I - |\psi\rangle\langle\psi|\}$ and any off-diagonal terms of the deviation do not contribute in the (survival) probabilities.

$$H_s = \sum_{ij} J_{ij} \left(\sigma_i^+ \sigma_j^- + \sigma_i^- \sigma_j^+ \right) + B \sum_j \sigma_j^z, \quad (9)$$

where $J_{ij} \propto \frac{J}{|i-j|^\alpha}$ is the interaction term and α dictates the strength of the interaction, $\sigma^{\pm, z}$ are the corresponding Pauli operators, N is the length of the spin chain, i.e. the number of qubits in the system, and B is the transverse magnetic field strength. We focus on two regimes; on the larger α values, where we assume $\alpha \rightarrow \infty$, i.e. nearest-neighbour interactions only [52], and on $\alpha \sim 0$ which corresponds to an all-to-all coupling between sites. Current experimental development [8, 53–55] has led to regimes where the theoretical prediction for the above XY Hamiltonian becomes classically intractable, it is therefore essential to provide a characterisation of such devices in the absence of simulation capabilities.

In the following we consider the case of a system consisting of $N = 6$ spins studying how the coupling strength, the disorder terms and the field strength affect the protocol. We simulate the evolution of the system from an initial product state $|\psi\rangle = |\uparrow\downarrow\uparrow\downarrow\downarrow\uparrow\rangle$, usually denoted as the charge-density wave state. We generate our set of unitaries $\{U_k = e^{-iH_k dt}\}$ by adding a disorder term chosen from equation (4) with $k = \{1, \dots, 1000\}$. In the first case (section 4.1), we run the ARB protocol with perfect inversion operators, and model gate and time-independent noise, adhering to the specifications of the standard RB protocol. We look at the impact of the field-term on the protocol and in the following section 4.2, we explore different noise scenarios such as *spontaneous emission*, *weakly time-dependent noise* and the case where we no longer have perfect inversion operators, and instead model the same type of noise in the backward evolution; as well as the impact of the time-step on our simplest case models.

4.1. Standard ARB with gate and time-independent noise

We modelled the gate and time-independent noise in the form of uniformly distributed fluctuations to both J and B terms that form the static Hamiltonian H_s (see equation (9)). These terms arise from fluctuations or calibration errors in the trapping fields and are native to the experimental device we model. Furthermore, we run the protocol with no errors in state preparation or measurement. Assuming that we have a sufficiently good ϵ -approximate 2-design, it follows that the decay curve that we expect our data to fit is of the form:

$$\begin{aligned} P_T &= A + Bf^T \\ &= \frac{1}{d} + \frac{d-1}{d}f^T, \end{aligned} \quad (10)$$

where $d = 2^N = 2^6$. In the following results, we fit the ARB decay curve using a non-linear least squares regression fitting tool [56]. Assuming no SPAM errors and perfect inverses simplifies the comparison of our results to those of an exact unitary 2-design (see theorem 3.2). We use this to estimate the average error rate by fitting P_T from the fidelity decay parameter f as $r = (d-1)(1-f)/d$.

4.1.1. XY Hamiltonian with transverse field (nearest-neighbours)

Here, we present the ARB fits for the model described in equation (9) for the case of $\alpha \rightarrow \infty$ that we consider to be well-described by:

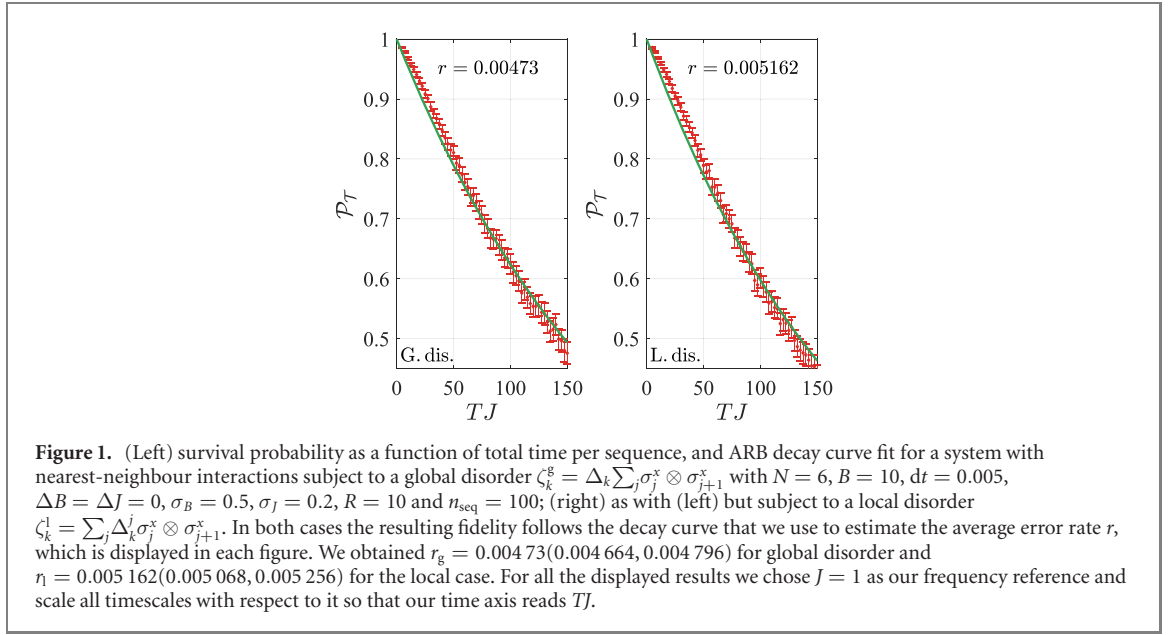
$$H_s = \sum_j J \left(\sigma_j^+ \sigma_{j+1}^- + \sigma_j^- \sigma_{j+1}^+ \right) + B \sum_j \sigma_j^z. \quad (11)$$

We generate a set of Hamiltonians $\{H_k\}$ for both local (ζ_k^l) and global (ζ_k^g) disorder terms:

$$\begin{aligned} H_k^g &= H_s + \Delta_k \sum_j \sigma_j^x \otimes \sigma_{j+1}^x \\ H_k^l &= H_s + \sum_j \Delta_k^j \sigma_j^x \otimes \sigma_{j+1}^x, \end{aligned} \quad (12)$$

where we chose from equation (4) that $u = x$, since this type of coupling should break the symmetry in our H_s , equation (11). This results in a set of $K = 1000$ unitaries (with $\Delta_k^{(j)}$ uniformly distributed with standard deviation $\delta = J$) from which to sample for each set $\{H_k^{(g,l)}\}$ of the following form:

$$\left\{ U_k^{(g,l)} = e^{-iH_k^{(g,l)} dt} \right\}. \quad (13)$$



In figure 1, we present the ARB protocol results for both the constant global disorder H_k^g (left) and the local site-dependent disorder H_k^l (right) as indicated in equation (12). On the time evolution forward the system is subject to noise proportional to H_s , which we chose to be normally distributed with mean $\Delta J = \Delta B = 0$ and standard deviations $\sigma_J = 0.2$ and $\sigma_B = 0.5^7$. In these results we conducted $n_{\text{seq}} = 100$ sequence iterations for every sequence time-length S_T and repeated each individual sequence $R = 10$ times to find the average of a given sequence. We discuss our choices for these parameters in appendix D.

In both cases, the data fits the ARB curve, though at earlier sequence time-lengths the data sits above the curve. The data fitting the curve could imply that the errors are depolarised by the process, which is conditioned on the set being a sufficiently good approximate 2-design. The fact that the data seems to fit to the curve at later times could indicate that the set of unitaries, both global and local, $\{U_k^{(g,l)}\}$ converge to a 2-design at these longer sequences. However, the fit of just one noise model to the curve is not sufficient to imply that our unitaries converge to a 2-design, which is why we explore more complex noise models in section 4.2. In fact, due to our noise model and running perfect inverses, it is possible that the averaging during ARB rather than twirling over an approximate 2-design is what causes the errors to behave like a depolarising channel; hence, we write our results in terms of theorem 3.2.

For the average error-rate we obtained (with 95% confidence bounds):

$$r_l = 0.005162(0.005068, 0.005256) \quad (14)$$

$$r_g = 0.00473(0.004664, 0.004796), \quad (15)$$

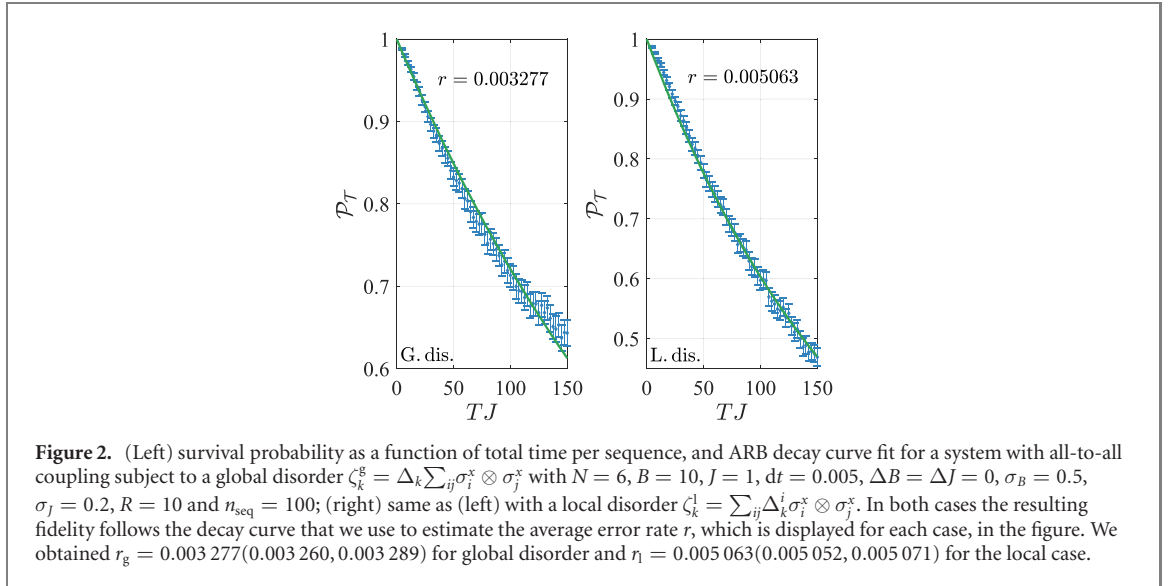
for locally (r_l) and globally (r_g) disordered unitaries, respectively. Where, from lemma F.5 we have the average error-rates bounded as follows:

$$(r_\mu - \epsilon) \leq r_\alpha \leq (r_\mu + \epsilon), \quad (16)$$

where $\alpha \in \{l, g\}$. In this case, r can be thought of as an average infidelity (due to no SPAM errors and no errors in the inversion operators [57]) per unit time when running a long sequence of these unitaries; therefore, a value of $r \ll 1$ indicates high average fidelity of running the unitaries intended. Since this type of benchmarking has not been applied in the analogue setting before, it is not clear what we would expect the average error rate to be. That is why we set a bound on our results (from theorem 3.2) and, in section 4.1.2 compare the r value obtained for one of our sets with the error per single unitary for different sets of initial states. It is reasonable to expect a small value of r with the noise model that we chose, and in the context of results gained from more refined noise models (section 4.2) our values of r_l and r_g here seem sensible.

As ARB has the potential to be implementable, with some modifications, it is relevant to discuss what r means if we do indeed have a sufficiently good approximate 2-design. We use r to characterise how this type

⁷ Note that these values are taken as an example and we do not require the experimental device to exhibit similar values.



of noisy hardware (modelled as the fluctuations to J and B) would behave under this set of unitary operators. In this case, the low values of $r_{(l,g)}$ obtained for our nearest-neighbour model would indicate that both of these sets of operations would perform well overall on such a hardware.

4.1.2. All-to-all spin model Hamiltonian with transverse field

Here, we present the ARB fits for a spin system governed by equation (9) for the case of $\alpha \sim 0$, i.e. an all-to-all coupled spin system with the same system parameters as in section 4.1.1:

$$H_s = \sum_{ij}^N J_{ij} (\sigma_i^+ \sigma_j^- + \sigma_i^- \sigma_j^+) + B \sum_j^N \sigma_j^z. \quad (17)$$

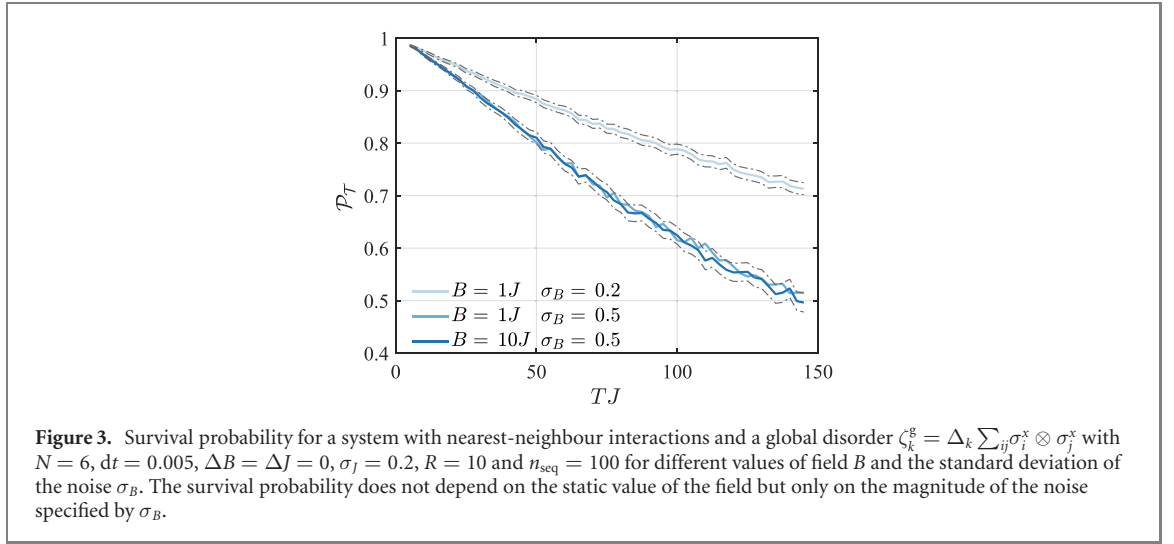
Again, we generated a set of Hamiltonians H_k for both local and global disorder terms:

$$\begin{aligned} H_k^g &= H_s + \Delta_k \sum_{ij} \sigma_i^x \otimes \sigma_j^x \\ H_k^l &= H_s + \sum_{ij} \Delta_k^{ij} \sigma_i^x \otimes \sigma_j^x, \end{aligned} \quad (18)$$

and thereby a set of 1000 unitaries $\{U_k^{(g,l)}\}$ for each $\{H_k^{(g,l)}\}$.

In figure 2 we present the results of fitting our data for the all-to-all $\{H_k^{(g,l)}\}$ to the ARB curve as in equation (10) with the globally disordered unitaries (left) and locally disordered unitaries (right). We observe that the two sets display substantially different profiles and, unsurprisingly, values of r . In the case of the globally disordered set, the data seems to fit the curve well at very early times then veering further away from the curve until it differs significantly at large TJ . This is in contrast to the profile of the locally disordered set, where we observe a much better agreement to the ARB curve. Moreover, the fit of the locally disordered all-to-all set is significantly better than the fit of either of the nearest-neighbour sets (figure 1), especially when comparing the behaviour at early times.

There are two results to address here, (i) comparison of the behaviour of the all-to-all sets $\{U_k^{(g,l)}\}$ and (ii) comparison of the behaviour of the locally disordered all-to-all set with both of the nearest-neighbour sets. Before discussing our thoughts on this, we again make it clear that the fit of *one* noise model on these sets does not offer a concrete conclusion as to whether we have an ϵ -approximate 2-design or not. However, if we assume that the data fitting the curve does imply convergence to a 2-design, due to this being a condition of RB producing meaningful results, then we can offer reasons for these results. To address (i) let us again consider the scrambling argument discussed in section 3.1. A global disorder term ζ_k^g applied to an all-to-all Hamiltonian of the form H_s (equation (17)) will not necessarily produce enough mixing in the Hamiltonians to explore the Hilbert space sufficiently. The good fit at early times is not enough to indicate that the protocol has depolarised the error channel, and could be explained as the cumulative error channel on much shorter time-sequences being closer to a depolarising channel when averaged, regardless of the structure of the gates applied.



To address (ii), as discussed earlier we would expect that the more random the disorder we create in our $\{H_k\}$ the faster the convergence to a 2-design with a large set of unitaries $\{U_k\}$, due to [34, 58]. The average error rates are as follows:

$$r_1 = 0.005\,063\,(0.005\,052, 0.005\,071) \quad (19)$$

$$r_g = 0.003\,277\,(0.003\,260, 0.003\,289), \quad (20)$$

with r_α bounded by equation (16). Again, if we assume a sufficiently good ϵ -approximate 2-design [an assumption partially substantiated in (right) figure 2] our average error-rates are low. Since we observed the best fit to the curve in the locally disordered all-to-all Hamiltonians, we wanted to determine how sensible this value of r_1 is. We numerically determined the average gate infidelity of the set of $K = 1000$ unitaries with this noise model, for two different sets of random initial states: N random product states and N random pure states, with $N = 100$. Since we assume that our set of unitaries converges to a 2-design after some time, it is not expected that the average error rate will be exactly the same; however, for random product states we found a value of $r_1 = 0.001\,50(0.001\,45, 0.000\,0155)$, which should be representative of the error for the state at the early stages, and for random pure states $r_1 = 0.010\,46(0.010\,38, 0.010\,54)$. We expect the latter to be higher than the actual value since a random pure state would exhibit substantially more entanglement than a state throughout the time evolution in our protocol. The fact that our value $r_1 = 0.005\,063$ sits between these two values demonstrates that it is a sensible estimate.

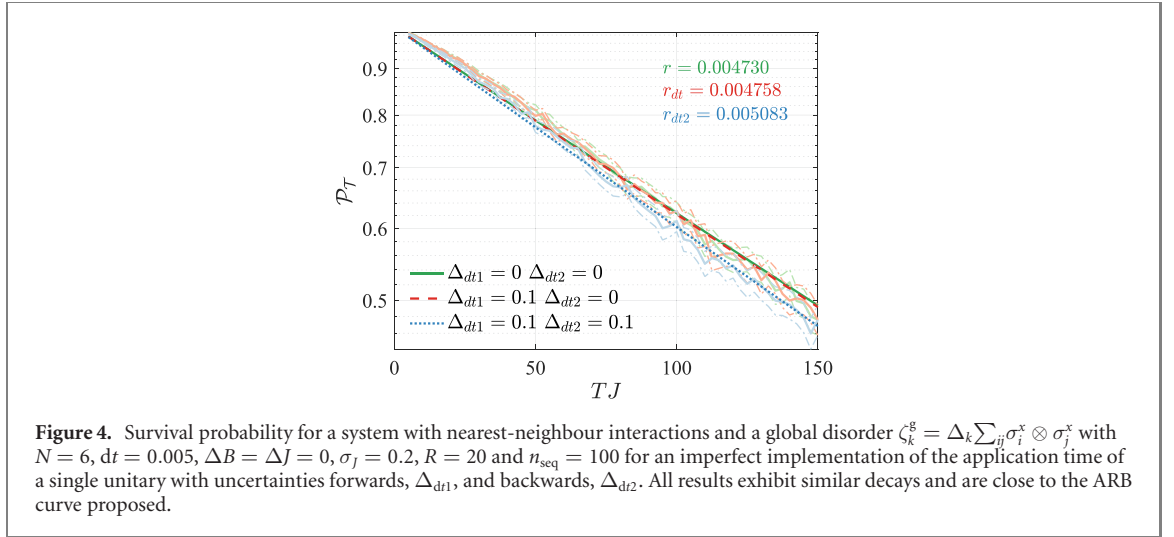
4.1.3. Impact of field-term

Obtaining the average error-rates (r) from the ARB protocol provides a characterisation of the hardware, and how it copes with a specific set of operations (gate-set). We therefore analyse whether some of the physical system parameters, of our specific tested system, may impact the measured values of r .

In figure 3, we present the fidelity decay curves for the case of the nearest-neighbour H_{XY} and added global disorder (see equations (11) and (12)) as a function of the transverse magnetic field, B . We observe that the ARB result does not depend on the off-set value (B) of the field but rather only depends on the magnitude of the noise specified by σ_B . This reveals that, according to our simulations, a quantity that would govern the ground state properties of the device does not affect our protocol. Therefore, the characterisation of the device depends only on the form of the noise and not on the choice of static parameters.

4.2. Further noise models

In search of more robust claims that our random sets $\{U_k\}$ approximate a 2-design and that ARB can provide meaningful results for more complex noise, we have two main considerations: (i) the survival probabilities revealed a better fit to the ARB curve after a finite time (apart from in the globally disordered all-to-all case) which could suggest that the scrambling over time causes the sets to converge to a 2-design. (ii) The copious averaging and simple noise model could be the reason for our errors presenting as a depolarising channel, and fitting the ARB curve. We investigate the latter by analysing the protocol in the presence of more elaborate noise: weakly time-dependent noise, spontaneous emission and imperfect inversion operators. These noise models are physically motivated and closer to the experimental conditions



of the implementation of the protocol. We also discuss the effects these noise models have on the decay curve in the context of convergence to a 2-design. For these analyses we chose the case of nearest-neighbour coupling with a global disorder term.

4.2.1. Weakly time-dependent noise

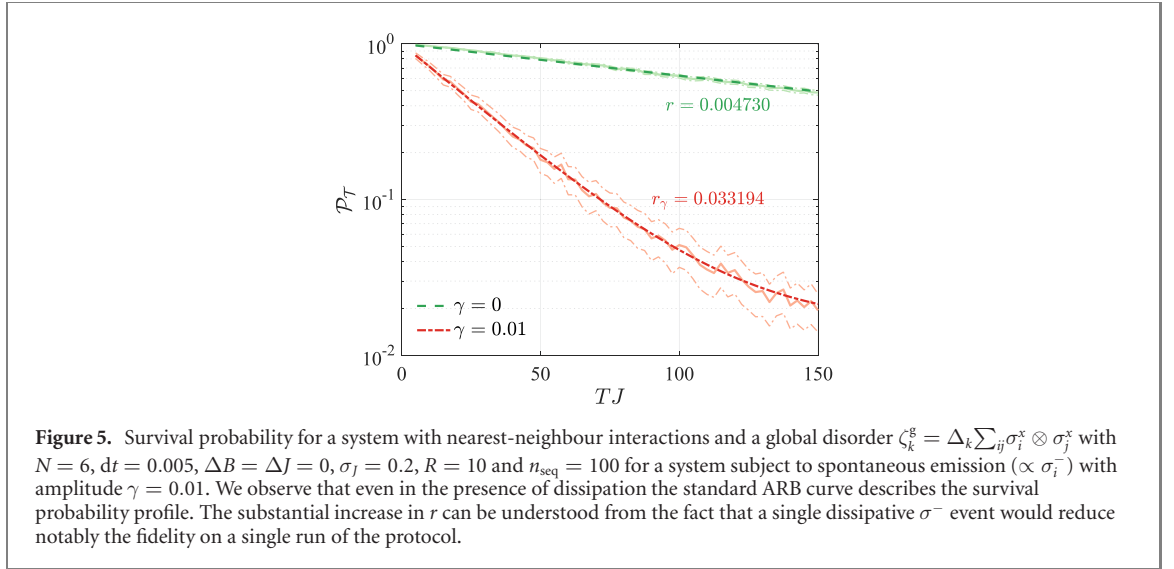
First, in figure 4 we study how robust our characterization is in the presence of some uncertainty in the time dt for which these random Hamiltonians are evolved and applied on the state. This could arise physically from a limited time resolution in the quantum hardware. We consider the case where every unitary is applied for a given time with an uncertainty of $\Delta_{dt1}/dt = 0.1$ in the forward evolution whilst inverted with the exact same time-step. We then consider the case when the backwards evolution has the same uncertainty $\Delta_{dt2}/dt = 0.1$ and so every H_k is inverted for a slightly different time. We compare these weakly time-dependent models to the fixed dt case, and observe that in all cases the decay profiles remain similar and fit the proposed ARB curve (see equation (10)). In the case of $\Delta_{dt2} = 0$ (dashed line), since any variation in the forward evolution is matched by the perfect time inversion the results remain unchanged and the estimated average error rate r_{dt1} is in the same confidence interval as r from fixed dt . Only when this uncertainty between the forward and backwards evolution differs (dotted line) do we observe a drop in the overall survival probability as expected from the unmatched time evolutions in both directions. Despite the presence of this weak time dependence the data exhibits similar agreement with the ARB curve which is a positive indication towards the reliability of our protocol.

4.2.2. Spontaneous emission

As mentioned previously, the affect of the ARB process on one noise model is not enough to indicate that our unitaries depolarise that channel. We therefore consider (again, on our globally disordered nearest-neighbour model) in figure 5 noise from spontaneous emission, an example of coupling of the quantum device to its environment. We model the time evolution of the open quantum system through quantum trajectories [59], with a ratio of $\gamma/J = 0.01$, which is compatible with experiments. We expect the average error rate for this noise model to be notably higher, as the dissipation term (which is $\propto \sigma_i^-$) that represents an emission event would transfer the state of the system to states with a much smaller overlap with the initial state. However, since the noise modelled is still chosen to be independent and uniform we expect that the results would fit the ARB curve, equation (10), if all necessary conditions are met by our generated set. We see that the presence of dissipation in the form of spontaneous emission does not impact the profile of the ARB decay curve, and as expected the average error rate is much larger than that found with no dissipation. This is yet another case that can be seen to indicate that the error channel is being simplified and characterised by our protocol.

4.2.3. Noisy time-inversion

Now that we have analysed how the ARB protocol is affected by both weak time-dependent noise and dissipation, it is necessary to address one of the assumptions of the implementation. Namely, in the previous results and those from section 4.2 we model the systematic time inversion as perfect. This choice was motivated by the fact that eliminating the time-inversion step is being explored as an extension to the project, in order to make the protocol more experimentally implementable. Nevertheless, this simplification can be regarded as unphysical in the present protocol and we therefore analyse the prospect

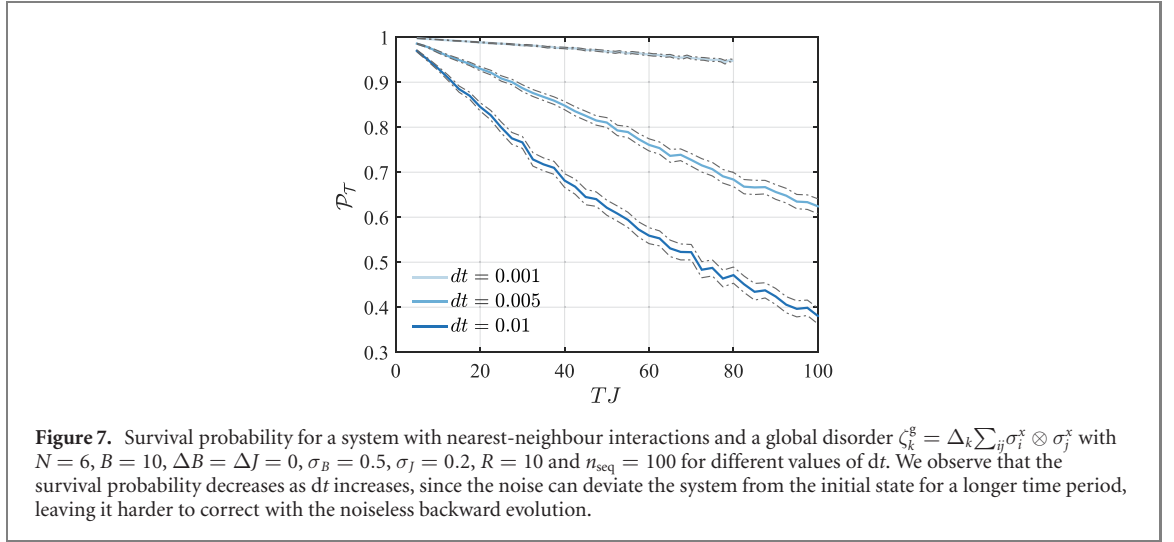
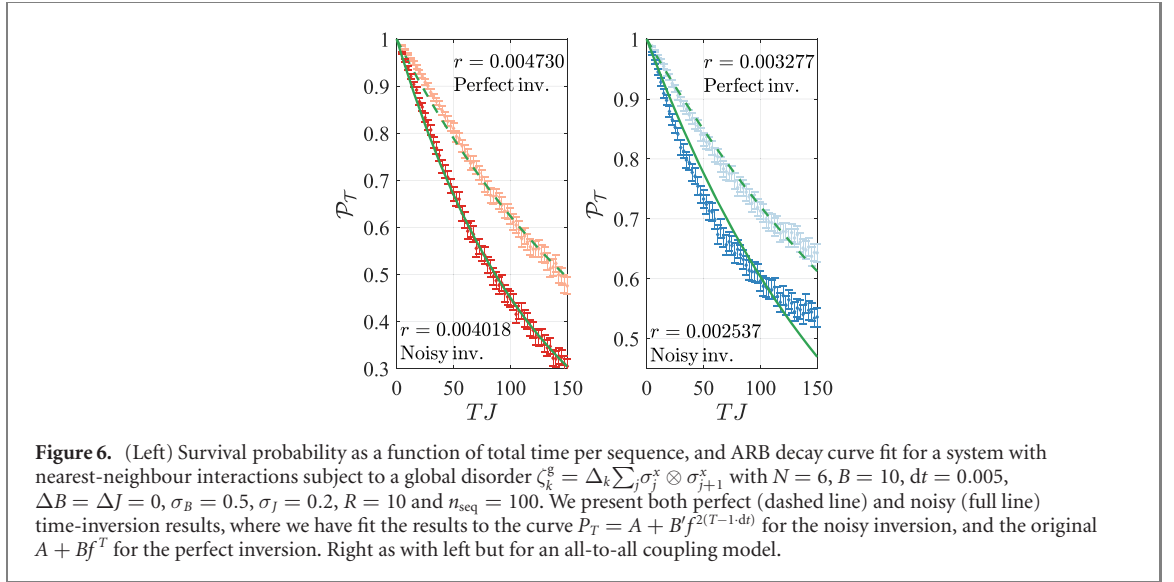


of a noisy time-inversion. We model the same type of noise (fluctuations to the J and B term) on the inversion operators. In this analysis, we consider both the nearest-neighbour and all-to-all coupling with the globally disordered set only, and compare the results to those in section 4.1. In figure 6 we compare the survival probability decay of both ideal and noisy time inversion for nearest-neighbour and all-to-all coupling. We fit the noisy inversion results to the following decay curve:

$$P_T = A + B' f_2^{2(T-1 \cdot dt)}, \quad (21)$$

where we set $B' = \frac{d-1}{d}$, following some reasonable assumptions, see appendix B. We now have noise in the backward evolution, as well as the forward, therefore adding an additional error channel after each inversion operator. Following the same proof as in the simpler case, with forward noise only, we end up with $(T - 1 \cdot dt)$ twirls of two error channels. This means that we have $(T - 1 \cdot dt)$ depolarising channels of the same strength, corresponding to the characterisation of a pair of error channels sequentially applied. We then get additional channels that do not scale with the time of the sequence T (see details in appendix B). In figure 6 and equation (21) we add a factor of two to the curve as we assume that the two composed error channels, when twirled, become a depolarising channel that is equivalent to the product of those two individually depolarised error channels. This is not an unrealistic approximation, and was made by Emerson et al [13] for a more restricted noise model. Interestingly, in the case of nearest-neighbours (left, figure 6) the results for noisy time inversion fit this curve much better than the perfect inversion fit its decay curve, which could be due to the factoring out of the SPAM-like errors that creates a closer fit at shorter sequences. We also observe that the average error rate $r_g = 0.004\,018(0.003\,982, 0.004\,054)$ is roughly 15%–16% smaller than the average error rate from the perfect inversion fit. This is not a huge difference, considering we assume an approximation that was developed for more limited noise, and that we make some approximations to B' . In this case, it seems that twirling the combined error channels does not result in the product of those two twirled error channels. Moreover, strictly speaking, the value we get as r_g can not be directly interpreted as the average gate infidelity since we now have added depolarising channels, and what we actually get here is the average error strength of two error channels (halved).

Comparing the noisy inversion results for the globally disordered all-to-all unitaries (right, figure 6), we see that the results veer even further from our derived decay curve. This is not surprising since we theorised that the unitary set for this model did not seem to converge to a 2-design, and twirling over the inverses would not be likely to change this. This noisy decay therefore gives some indication that our theory was correct, since it differs substantially from the ARB curve, especially at longer times. A final point that is worth mentioning here, is that in the more general and experimentally relevant case, where the inversion error is different than the forward error $\Lambda_{e'} \neq \Lambda_e$, our analysis would again provide an average error rate of a forward and backward evolution $(\Lambda_{e'} \circ \Lambda_e)$. This can then be used in different ways. For example, we could assume similar strength of noise in both directions and divide by two, and realistic implementations would have that $\Lambda_{e'} \sim \Lambda_e$ therefore this method would give us a reasonable average error rate. We could also use this to get an upper bound on the noise of the forward evolution, if we were to assume that all the errors came from the forward evolution.



4.2.4. Impact of time-step

Finally, we consider how the numerical time-step (dt) can impact the ARB curve. In figure 7, we present the survival probability \mathcal{P}_T for different values of the time-step (dt) used to create the unitaries for the same system as in section 4.1.1; again the nearest-neighbour H_{XY} with global disorder. The values of dt were chosen in the regime where the numerical simulations have no impact on the differences between them if the system were noiseless, to avoid any numerical error contribution to the analysis. These results highlight the fact that when a given noise term is applied for a longer period of time in the system it can cause a stronger deviation from the initial state, more difficult to correct with the perfect backwards evolution. In this analysis dt is, therefore, related to the ratio of change of the noise in the actual quantum device, which can affect the result since we model perfect inverses, i.e. noiseless backward evolution. For the purposes of the protocol, that the error on each unitary (gate) should not be dependent on the time it takes to run that unitary, we chose one value of $dt (= 0.005)$ fixed for all time-evolution's simulated.

4.3. Results overview

The results in section 4.2 support the conclusions that were drawn from the simplest noise model, section 4.1. Having observed that the unitary sets in the nearest-neighbour example seemed to converge to an approximate 2-design after some time, t , we found that the ARB protocol was robust to more complex noise for the globally disordered nearest-neighbour unitaries which reinforces the notion that this set converges to a 2-design. Furthermore, our simulations in section 4.1 reveal that the value of the magnetic field B that governs the state properties of the Hamiltonian, does not affect the results of our protocol; only the form of the noise added to B creates an affect. This supports the idea that our protocol is providing a measure of the noise in our (simulated) system and is robust to changes in the system parameters. The

fact that the protocol with these types of unitary sets is robust to weak-time dependent noise, and dissipation is important for experimental implementation. Additionally, we notice that in the case of the globally disordered all-to-all unitaries, which we theorise in section 4.1 do not converge to an approximate 2-design, the noisy inversion operators only bring the data further away from the ARB curve. This could indicate that the failure signatures are stronger when more complex noise, that does not necessarily adhere to the standard noise assumptions, is tested. Ultimately, we have created a version of randomized benchmarking with the physical capabilities of analogue quantum simulators in mind, and have found that the adapted protocol displays the behaviour we would expect under the noise tested and supports the notion that some of our sets approximate a 2-design.

5. Conclusions and future work

With the aim of developing a scalable generic method for testing analogue quantum simulators, we extended randomized benchmarking to the analogue setting. By replacing the quantum logic gates in the protocol with unitary time-evolution operators (native to the quantum system) requiring that they converge to a unitary 2-design, fixing the time-step to be the same for each unitary, and systematically inverting the unitaries rather than applying one single inversion operator, we presented the analogue randomized benchmarking protocol. In the context of continuous time evolution, the challenges we met were in: (i) creating a set of unitaries $\{U_k\}$ that generated an ϵ -approximate 2-design and understanding how the convergence rate of these unitaries affects the protocol; (ii) the generation of an efficient time-reversal of the unitaries on an analogue system. We numerically simulated our protocol on two models of the XY Hamiltonian (nearest-neighbour and all-to-all), which is native to trapped ions, adding both global and local disorder to generate the random unitary sets. We first modelled uniformly distributed fluctuations (noise) in the coupling J and B field terms of the static Hamiltonian. For the nearest-neighbour sets, the results fit the derived (for this noise) randomized benchmarking fidelity decay curve, particularly in the case of global disorder; this in turn indicated that the sets approximated a 2-design. For the globally disordered all-to-all case, the results did not fit the curve and it seemed that this set was not converging to a 2-design. We found the best fit to the decay curve for the local disorder (all-to-all), which we theorise is due to the richer dynamics of this set and supports the notion that this set converges faster to a 2-design. Therefore, in this case, we compared the average error rate to the average infidelity per gate and found the error rate predicted by our protocol was as expected for this type of gate and time-independent noise. Moreover, the robustness of the ARB decay curve was tested against weakly time-dependent noise, dissipation and an imperfect time-reversal scheme. We observed that for all scenarios of the globally disordered nearest-neighbour set, the proposed decay was suitable to describe the noise channels. For the globally disordered all-to-all unitaries, the imperfect time-reversal revealed further deviation from the curve, giving credence to our interpretation that this set does not converge to a 2-design.

Analogue randomized benchmarking creates opportunities for improving confidence in analogue quantum simulators by providing alternatives to the current benchmarking techniques. Assuming one has an ϵ -approximate 2-design and the sequences can be efficiently inverted, we could compare the average error-rate r across two quantum devices with the same starting Hamiltonian (H_s) that the set is built around; since ARB is primarily a test of a specific quantum hardware, r could provide information about what kind of noise were present in each device depending on the results of the protocol on both. Another area that ARB could be useful in is random circuit sampling, where the ARB parameter r could potentially be used to prove that random sampling from a random circuit is hard; with future works looking at this direction. At the root, ARB provides a measure for the performance of a set of unitaries on a specific hardware, and in the analogue setting this could be useful in testing programmable analogue quantum simulators. Particularly, the value of r would give a characterisation of how ones device will run a family of Hamiltonians, providing an extra security in the results you would gain from a programmable AQS experiment.

Extending RB to the analogue setting highlighted many interesting research questions, particularly in regards to approximate unitary t -designs with unitary time-evolution operators. In our work, we assume an ϵ -approximate 2-design is formed from our disordered set of unitaries $\{U_k\}$ (formed from disordered Hamiltonians $\{H_k\}$) because the disorder added was such that it should be sufficient to break the symmetry of the system Hamiltonian. However, we have not formally proven that our unitary sets $\{U_k\}$ are ϵ -approximate 2-designs and we therefore introduced a bound on the results garnered from the ARB protocol. This at least allows us to assess our results for the average error rate within a relative context, and with the standard error on our result we bound the unknown parameter ϵ . Perhaps the RB parameter r could provide an indication of the value of ϵ for a set of unitaries that categorically are an ϵ -approximate 2-design. An extension to this work is to formally define generating an ϵ -approximate unitary 2-design from

a set of unitaries formed around a Hamiltonian native to an AQS. The relations between the frame potential (see equation (E2)) and the Haar moment operator (see equation (E4)) [60, 61] that more accurately characterises an ϵ -approximate 2-design could provide a way to optimise the generation of approximate designs in the analogue setting. Moreover, exploring the types of disorder that one can add to the starting Hamiltonian, i.e. more locally-addressed, could reveal the optimal type of disorder that generates an ϵ -approximate 2-design with a given Hamiltonian. Furthermore, in this work we have shown that we can provide reasonable fits to the survival probability decay in the analogue setting for small system sizes; however, a relevant point to investigate is how the sampling, both in number of sequences n_{seq} and repetitions per sequence R , would exactly scale as a function of the system size.

Another area to investigate is the limitation of the time-reversal (mentioned in section 3.2) where we have acknowledged that systematic inversion could still provide a measure of the average error and that the main obstacle, in our point of view, to implementing our protocol is the fact that time-reversal in analogue devices is currently not feasible, although it can be implemented for a restricted set of operators, e.g. field terms. For small scale systems, one can compute the ideal output of running sequences on that system and estimate the fidelity of the output state with the ideal state, i.e. using DFE techniques or efficient tomography. This would mitigate the need for the inversion step in our protocol, and the benefit with this type of hybrid technique would be removing the SPAM errors from the characterization; though, unfortunately, losing the scalability advantage of ARB. On trapped-ion simulators in particular, digital and analogue computations may be performed, and therefore it would be prudent to look at the difference in errors found with both techniques: a possibility for ARB would be to implement the inversion in a trotterised (digital) way and combine the analogue and digital techniques in order to better characterise the types of errors on this kind of device.

The advantages that DRB (section 2) and ARB (section 3) have in common are that they evaluate, in a scalable way, the performance of a device whilst also removing the fixed imperfection of the SPAM errors. Comparatively, the use of native gates of the systems means that it is likely that ARB will have smaller errors, e.g. in compilation of gates/more noise that does not adhere to RB assumptions, than digital RB. This could be especially prevalent when dealing with the same physical system used for both analogue and digital quantum simulations. The assumption of gate-independence, and even of nearly gate-independence, of the noise model is far better motivated (and closer to reality) in the analogue case which means it is more likely that when experimentally implemented, ARB would give a better fit to the fidelity curve than in the digital case. Moreover, ARB could provide a way to test the performance of digital quantum simulators where researchers could focus on the average error-rate per length of computation time, rather than per-gate. This type of characterization is not only more physically motivated, but could also bring this analysis closer to the adiabatic model of quantum computation, where complexity is considered in regards to the time taken for the adiabatic evolution.

Acknowledgments

We thank the anonymous reviewers for their thoughtful comments and time. We thank Ulysse Chabaud, Andreas Elben, Martin Kliesch, Rawad Mezher and Hendrik Waldner for helpful discussions and clarifications. Work at the University of Strathclyde was supported by the EPSRC through the Programme Grant DesOEQ (EP/P009565/1) and through the UK Hub in Quantum Computing and Simulation (EP/T001062/1) and by the European Union's Horizon 2020 research and innovation programme under grant agreement No. 817482 PASQuanS. E K acknowledges support from the following: EPSRC Verification of Quantum Technology grant (EP/N003829/1) and UK Quantum Technology Hub: NQIT grant (EP/M013243/1) and the EU Flagship Quantum Internet Alliance (QIA) project. E D acknowledges support from the Doctoral Training Partnership (EP/N509711/1) under project No.1951737.

Appendix A. Randomized benchmarking

Running a unitary gate U on a physical device corresponds to a quantum channel denoted as Λ_U . The action of this quantum channel can be decomposed in two parts:

$$\Lambda_U := \Lambda_{U,e} \circ U, \quad (\text{A1})$$

where $\Lambda_{U,e} = \Lambda_U \circ U^\dagger$, capturing the errors that differentiate U from Λ_U . We use this convention to describe the imperfect channel as one that firstly applies the correct gate (U) followed by $\Lambda_{U,e}$, the error superoperator. In our protocol, we define $\Lambda_{U_{k_i}}$ as the imperfect implementation of a chosen unitary U_{k_i} .

To introduce a term commonly used in the literature, the probability that an initial state survives a quantum process is known as the *survival probability*. The survival probability of a channel Λ_{U_C} , where U_C is any unitary circuit, given a fixed initial state ρ_ψ is:

$$P := \langle \psi | \Lambda_{U_C}(\rho_\psi) | \psi \rangle = \text{Tr} (E_\psi \Lambda_{U_C}(\rho_\psi)), \quad (\text{A2})$$

where E_ψ is the projection on the state $|\psi\rangle$. Specifically in RB, we apply a sequence of imperfect unitaries followed by their (imperfect) inverse so that we have:

$$P = \langle \psi | \Lambda_{U^\dagger, e} \circ U^\dagger \circ \Lambda_{U, e} \circ U(\rho_\psi) | \psi \rangle. \quad (\text{A3})$$

It is clear to see that this probability is equal to unity if both the noise of the forward $\Lambda_{U, e}$ and the backward channels $\Lambda_{U^\dagger, e}$ are the identity (noiseless), since in that case we just evolve the state ρ_ψ by $U^\dagger \circ U = \mathbb{I}$. The average fidelity of the quantum channel (over all pure states) is defined as:

$$F(\Lambda_U, U) = \int d\psi \langle \psi | U^\dagger \Lambda_U(|\psi\rangle\langle\psi|) U | \psi \rangle, \quad (\text{A4})$$

and the average fidelity of a gate-set is given by $\int_U d\mu(U) F(\Lambda_U, U)$. The relevant quantity that we are interested in extracting is the average error-rate of a gate-set (on a specific hardware) which is simply one minus the average fidelity of the gate-set:

$$r := 1 - \int_U d\mu(U) F(\Lambda_U, U). \quad (\text{A5})$$

Why it works:

Experimentally, we obtain the average survival probability P_l for each length l (step 3), summing over all the sequences of the same length l . By the 2-design property of our unitary set $\{U_k\}$ we have:

$$\begin{aligned} P_l &= \frac{1}{n_{\text{seq}} \text{ sequences}} \sum \langle \psi_0 | \Lambda_{U_{\text{tot}}^\dagger} \Lambda_{U_l} \circ \Lambda_{U_{l-1}} \circ \cdots \circ \Lambda_{U_1}(\rho_{\psi_0}) | \psi_0 \rangle \\ &= \int dU_1 \cdots \circ dU_l \langle \psi_0 | \Lambda_{U_{\text{tot}}^\dagger} \Lambda_{U_l} \circ \Lambda_{U_{l-1}} \circ \cdots \Lambda_{U_1}(\rho_{\psi_0}) | \psi_0 \rangle, \end{aligned} \quad (\text{A6})$$

where $|\psi_0\rangle$ is the initial state of the system. Note that we have expressed the imperfect inversion operator as a single gate $\Lambda_{U_{\text{tot}}^\dagger}$. Decomposing the errors and assuming that they are gate and time-independent $\Lambda_{U, e} = \Lambda_e$, leads to:

$$P_l = \int dU_1 \cdots dU_l \langle \psi_0 | \Lambda_s \circ U_1^\dagger \cdots \circ U_l^\dagger \circ \Lambda_e \circ U_l \circ \Lambda_e \circ U_{l-1} \cdots \circ \Lambda_e \circ U_1(\rho_{\psi_0}) | \psi_0 \rangle. \quad (\text{A7})$$

Integrating over U_l twirls one channel $\Lambda_e \rightarrow \Lambda_{e, t}$, where $\Lambda_{e, t}$ is the depolarised (twirled) channel corresponding to Λ_e and the probability that characterises this channel is p_e (see equations (A9) and (A11)). One can then integrate one-by-one the U_k 's, where each of the integrals result in one error term being twirled and hence depolarised. Noting that the twirled (depolarised) channels commute with all other channels in general and specifically with the unitaries appearing in the above expression, we obtain:

$$P_l = \langle \psi_0 | \Lambda_s \circ (\Lambda_{e, t})^l(\rho_{\psi_0}) | \psi_0 \rangle. \quad (\text{A8})$$

Here, Λ_s represents the error channel corresponding to the SPAM errors. Since the imperfect inverse $\Lambda_{U_{\text{tot}}^\dagger}$ is one single operator (or at the very least it will be composed of far less gates than the forward sequence) the error associated with it can be absorbed into the SPAM errors. These errors can also be treated as a depolarising channel, because the state is measured in the basis $\{|\psi\rangle\langle\psi|, I - |\psi\rangle\langle\psi|\}$ and the corresponding 'off-diagonal' terms do not affect the probabilities that we measure (and need for the subsequent estimations). This SPAM error depolarising channel (Λ_s) is characterised by the parameter p_s , leading to:

$$\begin{aligned} P_l &= p_s p_e^l + (1 - p_s p_e^l) \frac{1}{d} \\ &= \frac{1}{d} + \left(\frac{d-1}{d} \right) p_s p_e^l, \end{aligned} \quad (\text{A9})$$

which is in the exact form $P_l = A + Bf^l$ mentioned in step 4, where $f = p_e$, $A = 1/d$, and $B = \left(\frac{d-1}{d}\right) p_e$. By plotting P_l for different values of l we recover the value of p_e (f). Having obtained the depolarising probability of the error-channel, we can now look at the average fidelity of the gate-set:

$$\begin{aligned} \int_U d\mu(U) F(\Lambda_U, U) &= \int_U d\mu(U) F(\Lambda_{U,e}, I) \\ &= \int_U d\mu(U) F(\Lambda_e, I), \end{aligned} \quad (\text{A10})$$

and due to the left-invariance of the Haar measure, we have that $F(\Lambda_e, I) = F(\Lambda_{e,t}, I)$, i.e. the fidelity of any superoperator (Λ_e) with the identity (I) is equal to the fidelity of its exact Haar twirl ($\Lambda_{e,t}$) with the identity (I) [62]. Therefore, with the simplifying assumptions made, it is clear to see how the average fidelity is related to p_e , since

$$\begin{aligned} \int_U d\mu(U) F(\Lambda_U, U) &= \int_U d\mu(U) F(\Lambda_{e,t}, I) \\ &= \int_U d\mu(U) \left(p_e + \frac{1-p_e}{d} \right) \\ &= \frac{1}{d} + \left(\frac{d-1}{d} \right) p_e. \end{aligned} \quad (\text{A11})$$

Recalling that $r := 1 - F_{\text{ave}}$ (see equation (A5)) we get the expression of step 5 for the average error-rate of the gate-set: $r = (d-1)(1-p_e)/d$.

Appendix B. Noisy time-inversion ARB

Here, we analyse the decay curve in the presence of noisy inversion operators. The sequences that we apply in this scenario are of the following form:

$$\Lambda_e \circ U_1^{-1} \circ \Lambda_e \circ \cdots \circ \Lambda_e \circ U_{l-1}^{-1} \circ \Lambda_e \circ U_l^{-1} \circ \Lambda_e \circ U_l \circ \Lambda_e \circ U_{l-1} \circ \Lambda_e \circ \cdots \circ \Lambda_e \circ U_1, \quad (\text{B1})$$

where for the error channel on both forward and backwards evolution, we assume gate and time-independence, i.e. $\Lambda_{U,e} \equiv \Lambda_e$ and have decomposed the errors, as in equation (A7). Now, writing the survival probability for this sequence we have:

$$\begin{aligned} P_l &= \int dU_1 \cdots dU_l \langle \psi_0 | \Lambda_e \circ U_1^{-1} \circ \Lambda_e \circ \cdots \circ \Lambda_e \circ U_{l-1}^{-1} \circ \Lambda_e \circ U_l^{-1} \circ \Lambda_e \circ U_l \circ \\ &\quad \Lambda_e \circ U_{l-1} \circ \Lambda_e \circ \cdots \circ \Lambda_e \circ U_1 | \rho_\psi \rangle | \psi_0 \rangle. \end{aligned} \quad (\text{B2})$$

Integrating over U_l results in three depolarizing channels. Firstly, the error channel in the ‘middle’ of the sequence is twirled over all the unitaries in that space (over all random unitaries in n_{seq} applied here) such that:

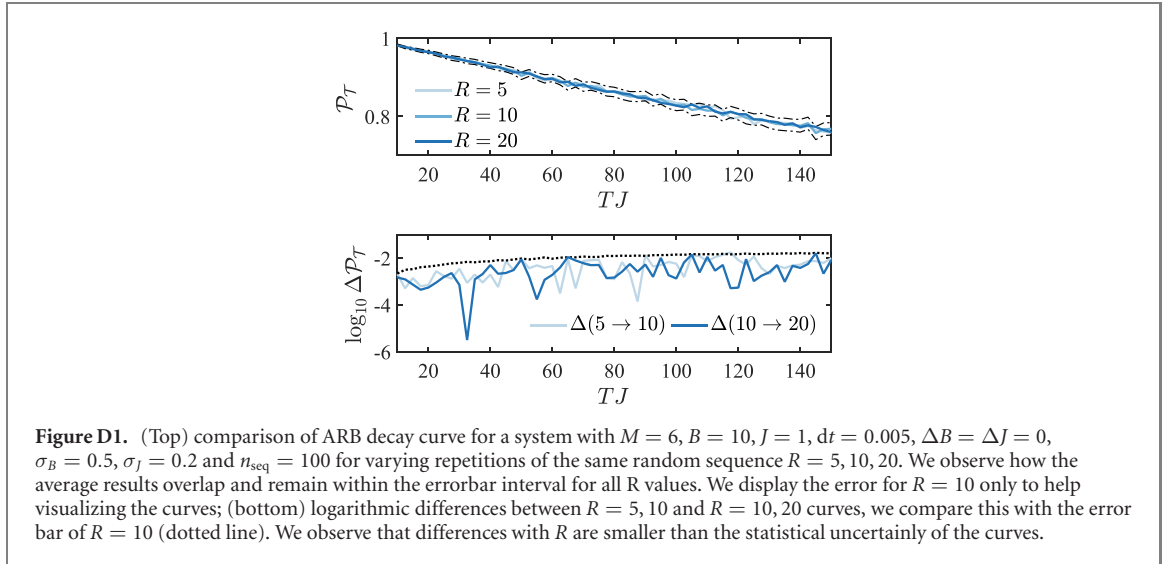
$$\int_U dU U_l^{-1} \circ \Lambda_e \circ U_l \rightarrow (\Lambda_e)_t = p_e + (1-p_e) \frac{1}{d}, \quad (\text{B3})$$

where $(\Lambda_e)_t$ is the depolarized (twirled) channel corresponding to Λ_e and the probability that characterizes that channel is p_e . For the rest of the twirls, it is two error channels that become depolarized, in the following way:

$$\int_U dU U_{l-1}^{-1} \circ \Lambda_e \circ \Lambda_e \circ U_{l-1}, \quad (\text{B4})$$

which results in $l-1$ depolarizing channels of the form $(\Lambda_e \circ \Lambda_e)_t = p_2 + (1-p_2) \frac{1}{d}$. The last depolarizing channel comes from the error channel composed on the first inversion operator (far left in equation (B1)). These errors can be seen as SPAM errors, and can be treated as a depolarizing channel in the same way, since they are independent of sequence length. We can denote the parameter that characterizes this depolarizing channel as p_e . The resulting survival probability for sequences of this form is:

$$\begin{aligned} P_l &= p_e p_2^{l-1} + (1-p_e p_2^{l-1}) \frac{1}{d} \\ &= \frac{1}{d} + \left(\frac{d-1}{d} \right) p_e p_2^{l-1}, \end{aligned} \quad (\text{B5})$$



which is in the form of equation (21), where $B' = \frac{d-1}{d} p_e p$ and $p_2 = f_2$. For simplicity and easier direct comparison with the perfect inversion, we will assume that $p = p_e = 1$. The actual values are very close to unity. We could have instead left these parameters as variables to be extracted from fitting the curve, something that would have a low impact on our results. The approximation we make in section 4.2.3 is that the twirl of two composed error channels is the same as the product of those error channels individually twirled (the square of the twirled error channel), i.e. $(\Lambda_e \circ \Lambda_e)_t = (\Lambda_e)_t (\Lambda_e)_t = (\Lambda_e)_t^2$. Therefore, in our fit we add the factor of two to the curve which gives a more direct comparison to the parameter r from the original curve.

Appendix C. Unitary 2-designs

Consider a superoperator Λ acting on a space M_D of D -dimensional quantum states, when $t = 2$, Λ has a $D^2 \times D^2$ dimensional matrix representation. We now define a set $U(D)$ of unitary matrices on this space M_D . If the set is a unitary 2-design then the space M_D is reducible to two irreducible invariant subspaces. Now, we define Λ acting on a quantum operator X as: $\Lambda(X) = AXB$. Schur's lemma [63] implies the following (reducible representation) for $U(D)$ -invariant trace-preserving operators:

$$\Lambda(X) = pX + (1 - p) \text{Tr}(X) \frac{\mathbb{I}}{D}, \quad (\text{C1})$$

where $p = \frac{\text{Tr}(\Lambda) - 1}{D^2 - 1}$. Considering the fact that a unitary 2-design means that sampling uniformly from the set $\{U_1, \dots, U_K\}$ is operationally equivalent to sampling from the Haar measure, we can say that [43]:

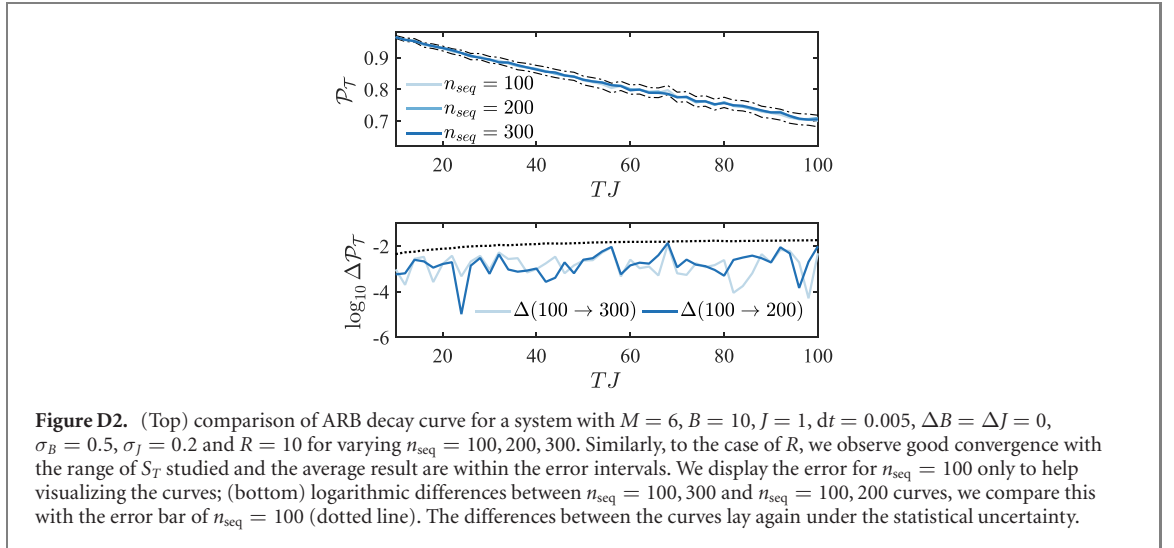
$$\frac{1}{K} \sum_{k=1}^K U_k^\dagger A U_k X U_k^\dagger B U_k = \int_{U(D)} dU U^\dagger A U X U^\dagger B U, \quad (\text{C2})$$

for all $A, X, B \in L(C^D)$. This essentially means that if we have a set $\{U_k\}$ that is a unitary 2-design or above, then conjugating a quantum channel over this set and averaging will result in a depolarisation of that channel.

Appendix D. Parameters convergence

Changes to some of the method parameters, such as the repeated runs of each random sequence (R) and the number of sequences tested for each sequence length (n_{seq}) would improve the accuracy of our results: iterations over as many sequences as possible of the same length are desired in RB to sample as uniformly from the unitary space as possible and repeating each sequence sufficiently gives a more accurate measure of the average survival probability for that sequence. In this section we discuss the choice of numerical parameters in the results presented in the main text. Here we describe how the RB curves depend on some of the averaging parameters such as R and n_{seq} .

In figures D1 and D2, we justify the choice of the numerical parameters $n_{\text{seq}} = 100$ and $R = 10$ in the main text by comparing the ARB decay curves for varying values of the mentioned variables. We observe



that in both cases the statistical uncertainty derived from the sequence averaging is larger than the discrepancy as we vary these parameters, therefore we are confident that the presented results do not depend on the chosen values for n_{seq} and R .

Appendix E. Comparative techniques for 2-designs

In addition to our discussion in section 3.1, here we highlight some of the comparative techniques used to determine whether one has an exact unitary 2-design beginning with the introduction of the spherical t -design.

Consider a real function f , and imagine we are interested in the average value of this function on an n -dimensional real sphere S^n ; this is hard to compute, so one can think of averaging over a finite set of unit vectors $D = \{|\phi_1\rangle, \dots, |\phi_K\rangle\}$ instead. Briefly, a spherical t -design is a finite subset D of S^n such that the average of every t th order polynomial p over S^n is equal to the average of p over D . For spherical t -designs, the frame potential [44] is a well-known metric for determining whether one has an exact spherical design or not, and is defined as follows:

Definition E.1 (Spherical t -design). A set of vectors $\{|\phi_1\rangle, \dots, |\phi_K\rangle\}$ is a spherical 2-design in \mathbb{C}^d if and only if:

$$\sum_{k,k'} \frac{|\langle \phi_k | \phi_{k'} \rangle|^4}{K^2} = \frac{2}{d^4 + d^2}. \quad (\text{E1})$$

The definition of spherical t -designs was modified for the unitary setting, changing the real sphere S^n to a set of unitaries $U(D)$ and comparing with the Haar distribution, with the term *unitary t -design* coined by Dankert [64].

Adapting the frame potential from the spherical setting to the unitary setting, Gross *et al* [65] developed the frame potential technique for determining whether ones set of unitaries is an exact unitary 2-design or not, see definition E.2. Another technique for determining whether the (suspected) unitary 2-design that one has in an exact unitary 2-design, is to compare the moment operators up to order t (in this case 2) of the particular unitary set and the Haar measure, we refer to [45] for definition E.3.

Definition E.2 (frame potential). Let the set $\mathbb{M} = \{U_k\}$ with $\{k = 1, \dots, K\}$ be a set of unitaries. The frame potential of \mathbb{M} is defined as:

$$P(\mathbb{M}) = \sum_{U_k, U_{k'} \in \mathbb{M}} \frac{|\text{Tr}(U_k^\dagger U_{k'})|^4}{K^2}, \quad (\text{E2})$$

\mathbb{M} is an exact unitary 2-design $\Leftrightarrow P(\mathbb{M}) = 2$.

Definition E.3 (second order moment operators). A degree (t, t) -monomial in $C \in U((\mathbb{C}^d)^{\otimes n})$ is degree t in the entries of C and degree t in the entries of C^* . Setting $t = 2$, and collecting all these monomials into a single matrix of dimension d^{2n2} by defining $C^{\otimes 2,2} := C^{\otimes 2} \otimes C^{*\otimes 2}$, we state that α is an exact unitary 2-design if expectations of all $(2, 2)$ moments of α match those of the Haar measure:

$$G_\alpha^{(2)} = \mathbb{E}_{C \sim \alpha} [C^{\otimes 2} \otimes C^{*\otimes 2}]. \quad (\text{E3})$$

Therefore, μ is an exact unitary 2-design if and only if:

$$G_\alpha^{(2)} = G_\mu^{(2)}, \quad (\text{E4})$$

Where μ is the Haar distribution.

Appendix F. Proofs for section 3.1

F.1. Proof of theorem 3.1

In order to prove theorem 3.1 we first present the following definitions:

Definition F.1. The trace-norm of a quantum channel \mathbb{E} in terms of the input state density matrix ρ , that minimises the error probability on distinguishing between two quantum channels \mathbb{E}_1 and \mathbb{E}_2 , is defined as:

$$\|\mathbb{E}\|_1 := \max_{\rho} \|\mathbb{E}(\rho)\|_1, \quad (\text{F1})$$

where $\|\cdot\|_1$ is the trace-norm, i.e. $\|X\|_1 = \text{Tr} \sqrt{X^\dagger X}$.

Definition F.2. The diamond-norm distance written in terms of the trace-norm of a quantum channel \mathbb{E} is as follows:

$$\begin{aligned} \|\mathbb{E}\|_\diamond &= \|I \otimes \mathbb{E}\|_1 \\ &\geq \|\mathbb{E}\|_1. \end{aligned} \quad (\text{F2})$$

Definition F.3. The average survival probability for each sequence length found from RB is:

$$P_l^\mu = A + Bf^l, \quad (\text{F3})$$

with a pure input state ρ , and using an *exact* unitary 2-design, where the average error rate of the unitaries $\{U_\mu\}$ is $r = (d-1)(1-f)/d$ (see conversations surrounding, and including, equation (A11)).

Similarly, the average survival probability for each sequence length, measured for an unknown α -distribution of unitaries via RB, is:

$$P_l^\alpha := P_l^\mu \pm \delta P_l, \quad (\text{F4})$$

where the unitaries are assumed to be an ϵ -approximate 2-design and $P_l^\tau = \text{Tr}[E_\psi \mathbb{E}_\tau(\Lambda)^l(\rho_\psi)]$, $\tau \in \{\alpha, \mu\}$ represents the survival probabilities of input state ρ_ψ when the twirled quantum channels \mathbb{E}_τ are applied to it l times.

Now, we present the following lemma's and their proofs.

Lemma F.1. If $\mathbb{E}_\alpha(\Lambda)$ is the twirled channel of Λ over a set of unitaries $\{U_\alpha\}$ spread according to a probability distribution α and \mathbb{E}_μ is the Haar-twirl of that channel, then for an ϵ -approximate 2-design, it holds that:

$$\|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1 \leq \epsilon, \quad (\text{F5})$$

with $\mathbb{E}_\alpha(\Lambda)(\rho) = \int_U d\alpha(U) U \circ \Lambda \circ U^\dagger(\rho)$.

Proof for lemma F.1. From definitions 3.1 and F.2 we obtain:

$$\|\mathbb{E}_\alpha(\Lambda) - \mathbb{E}_\mu(\Lambda)\|_1 \leq \|\mathbb{E}_\alpha(\Lambda) - \mathbb{E}_\mu(\Lambda)\|_\diamond \leq \epsilon, \quad (\text{F6})$$

which implies:

$$\|\mathbb{E}_\alpha(\Lambda) - \mathbb{E}_\mu(\Lambda)\|_1 \leq \epsilon. \quad (\text{F7})$$

Using definition F.1 we state:

$$\|\mathbb{E}_\alpha(\Lambda) - \mathbb{E}_\mu(\Lambda)\|_1 := \max_{\rho} \|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1. \quad (\text{F8})$$

And it holds by definition that:

$$\|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1 \leq \max_{\rho} \|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1. \quad (\text{F9})$$

Therefore:

$$\|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1 \leq \max_{\rho} \|\mathbb{E}_\alpha(\Lambda)(\rho) - \mathbb{E}_\mu(\Lambda)(\rho)\|_1$$

$$\leq \|\mathbb{E}_\alpha(\Lambda) - \mathbb{E}_\mu(\Lambda)\|_1 \leq \epsilon, \quad (\text{F10})$$

as required. \square

Lemma F.2. *If the unitaries $\{U_\alpha\}$ form an ϵ -approximate 2-design, it holds that:*

$$|P_l^\alpha - P_l^\mu| = |\delta P_l| \leq l \cdot \epsilon. \quad (\text{F11})$$

Proof for lemma F.2. Considering a fixed length l , we define the state after a sequence of this length, S_l , of imperfect unitaries $\Lambda_{U_l} = \Lambda_\epsilon \circ U_l$, has been applied to initial state $\rho_\psi = |\psi\rangle\langle\psi|$ and before a measurement has been taken, as:

$$\rho(S_l) := \Lambda_\epsilon \circ U_1^\dagger \cdots U_l^\dagger \circ \Lambda_\epsilon \circ U_l \circ \Lambda_\epsilon \circ U_{l-1} \cdots \circ \Lambda_\epsilon \circ U_1(|\psi\rangle\langle\psi|). \quad (\text{F12})$$

We set ρ_0 as the average state of the RB protocol before the final measurement (i.e. averaging the above expression over different sequences according to the distribution $d\alpha$):

$$\rho_0 = \int d\alpha(U_1) \cdots d\alpha(U_l) \rho(S_l). \quad (\text{F13})$$

which is equivalent to the average channel defined in lemma F.1 but for a sequence of length l , i.e. $E_\alpha(\Lambda)^l(\rho)$. Similarly, we define the following:

$$\begin{aligned} \rho_1 &= \int d\alpha(U_1) \cdots d\alpha(U_{l-1}) d\mu(U_l) \rho(S_l) \\ &\cdots \\ \rho_j &= \int d\alpha(U_1) \cdots d\alpha(U_{j-1}) d\mu(U_j) \\ &\cdots d\mu(U_l) \rho(S_l) \\ &\cdots \\ \rho_l &= \int d\mu(U_1) \cdots d\mu(U_l) \rho(S_l). \end{aligned} \quad (\text{F14})$$

Here, in each of the above quantum states we replace (one-by-one) the average over the distribution α with that of the Haar measure μ . Any two consecutive states ρ_j, ρ_{j+1} differ by a single integration, and by the ϵ -approximate 2-design property (see equation (F5)) we, therefore, have that:

$$\|\rho_j - \rho_{j+1}\|_1 \leq \epsilon \implies \|\rho_l - \rho_0\|_1 \leq l \cdot \epsilon, \quad (\text{F15})$$

where the implication follows from the triangle inequality. From the definition of the trace-norm we get:

$$\begin{aligned} |\langle\psi|\rho_l|\psi\rangle - \langle\psi|\rho_0|\psi\rangle| &\leq l \cdot \epsilon \\ |P_l^\mu - P_l^\alpha| &\leq l \cdot \epsilon. \end{aligned} \quad (\text{F16})$$

The definition of survival probability stated in equation (A2) results in the the equivalence of the above equation. Intuitively, the difference in the probabilities, $\langle\psi|\rho_l|\psi\rangle$ and $\langle\psi|\rho_0|\psi\rangle$, of obtaining the states (after measurement) is no larger than $l \cdot \epsilon$, the bound specified by the trace-norm between the the two initial states, as in equation (F15). \square

Thus completing the proof.

F.2. Proof of theorem 3.2

Lemma F.3. *Let A and B be known quantities. Under assumption 3.1 and with a small ϵ , the error in determining f from the RB method of an ϵ -approximate 2-design, is given by:*

$$\delta f_l \approx \frac{\epsilon}{f^{l-1}B}, \quad (\text{F17})$$

where f is dependent on l , $f_l = f \pm \delta f_l$ is the value for the fidelity decay parameter found for RB with an ϵ -approximate 2-design and f is that obtained with an exact design.

Note that in the following proofs we denote the error in obtaining a quantity C as δC , where a subscript is added if the error depends on some measured quantity that is not obvious.

Proof for lemma F.3. From definition F.3 we have that:

$$P_l^\mu \pm \delta P_l = A + B(f \pm \delta f_l)^l. \quad (\text{F18})$$

With lemma F.2 and under assumption 3.1, it follows that:

$$\begin{aligned} P_l^\mu \pm l \cdot \epsilon &= A + B(f \pm \delta f_l)^l \\ f^l \pm \frac{l \cdot \epsilon}{B} &= (f \pm \delta f_l)^l \\ f^l \pm \frac{l \cdot \epsilon}{B} &\approx f^l \pm l \cdot f^{l-1} \delta f_l, \end{aligned} \quad (\text{F19})$$

where the last approximation holds if $\delta f_l/f \ll 1$. This leads to:

$$\delta f_l = \frac{\epsilon}{f^{l-1}B}. \quad (\text{F20})$$

□

Note that f_l and δf_l can be computed separately for each different length l . In practice, B will also depend on the SPAM errors which are (generally) unknown, and therefore it is essential to consider several different lengths to obtain a value for f_l .

Lemma F.4. *The error in determining the average error-rate r_l of a gate-set that is an ϵ -approximate 2-design, from the RB method, is given by:*

$$\delta r_l = \frac{\epsilon}{p_s f^{l-1}}, \quad (\text{F21})$$

as compared to that when using RB with an exact 2-design. Where p_s is the parameter that characterises the depolarising channel of the SPAM errors (i.e. if no SPAM errors are present, $p_s = 1$, while if SPAM errors completely depolarise the channel $p_s = 0$).

Proof for lemma F.4. Similarly, given that $r = \frac{d-1}{d}(1-f)$, $B = \frac{(d-1)p_s}{d}$ and the previous result $\delta f_l = \frac{\epsilon}{B f^{l-1}}$, we obtain:

$$\begin{aligned} r \pm \delta r_l &= \frac{B}{p_s}(1 - (f \pm \delta f_l)) \\ \delta r_l &= \frac{\epsilon}{f^{l-1}p_s}. \end{aligned} \quad (\text{F22})$$

□

It is clear to see that the larger the l that we use to estimate r_l , the larger the error in that estimation due to the gate-set being an ϵ -approximate 2-design. In our simulated results, we assume no SPAM errors ($p_s = 1$), and we can therefore take the more optimistic view and consider the errors for $l = 1$. Under these assumptions, it is easily seen that:

Lemma F.5. *If $l, p_s = 1$ the value of r determined from the RB method when testing an ϵ -approximate 2-design, is bounded as follows:*

$$\begin{aligned} r' &= r \pm \delta r \\ r - \epsilon &\leq r' \leq r + \epsilon, \end{aligned} \quad (\text{F23})$$

where we use r' to denote the value for the average error-rate measured from using an ϵ -approximate 2-design, and r is the value gained when using an exact 2-design, with the RB method.

For simplification, we assume no SPAM errors, $p_s = 1$, in our analysis; however, it is clear that the error in f_l and r_l would increase with greater SPAM errors. By extrapolating the values of f' and r' from the survival probability of the smallest length, $l = 1$, equation (F22) allows us to find the smallest error. In practise the SPAM errors p_s exist and we need multiple values of l to estimate and remove this contributing factor from p_s . This is not necessary for our purposes and we can take the weakest bound on the average error rate, where $l = 1$, and is given by $\delta r_{\min} = \epsilon$, while $\delta f_{\min} = \epsilon/B$.

Proof for lemma F.5. When we set $l = 1$ and $p_s = 1$, the error induced in the fidelity decay parameter f and the average error-rate r become:

$$\delta f = \frac{\epsilon}{B} \quad (\text{F24})$$

$$\delta r = \epsilon, \quad (\text{F25})$$

and therefore:

$$r - \epsilon \leq r' \leq r + \epsilon, \quad (\text{F26})$$

where $r' = r \pm \delta r$ is the average error-rate found from RB using an ϵ -approximate 2-design and r is that found from an exact 2-design. \square

ORCID iDs

E Derbyshire  <https://orcid.org/0000-0002-7901-0724>

P Wallden  <https://orcid.org/0000-0002-0255-6542>

References

- [1] Kokail C *et al* 2018 arXiv:1810.03421
- [2] Yang D, Grankin A, Sieberer L M, Vasilyev D V and Zoller P 2019 arXiv:1905.06444
- [3] Hangleiter D, Kliesch M, Schwarz M and Eisert J 2017 *Quantum Sci. Technol.* **2** 015004
- [4] Bermejo-Vega J, Hangleiter D, Schwarz M, Raussendorf R and Eisert J 2018 *Phys. Rev. X* **8** 021010
- [5] Cramer M, Plenio M B, Flammia S T, Gross D, Bartlett S D, Somma R, Landon-Cardinal O, Liu Y-K and Poulin D 2010 *Nat. Commun.* **1** 149
- [6] Mohseni M, Rezaekhani A T and Lidar D A 2008 *Phys. Rev. A* **77** 1094–622
- [7] Flammia S T and Liu Y-K 2011 *Phys. Rev. Lett.* **106** 230501
- [8] Lanyon B P *et al* 2016 *Nat. Phys.* **13** 1158–162
- [9] Torlai G, Mazzola G, Carrasquilla J, Troyer M, Melko R and Carleo G 2018 *Nat. Phys.* **14** 447
- [10] Toth G, Wieczorek W, Gross D, Krischek R, Schwemmer C and Weinfurter H 2010 *Phys. Rev. Lett.* **105** 250403
- [11] Gardiner S A, Cirac J I and Zoller P 1997 *Phys. Rev. Lett.* **79** 4790
- [12] Gorin T, Prosen T, Seligman T H and Žnidarič M 2006 *Phys. Rep.* **435** 33
- [13] Emerson J, Alicki R and Życzkowski K 2005 *J. Opt. B: Quantum Semiclass. Opt.* **7** S347–52
- [14] Knill E, Leibfried D, Reichle R, Britton J, Blakestad R B, Jost J D, Langer C, Ozeri R, Seidelin S and Wineland D J 2008 *Phys. Rev. A* **77** 012307
- [15] Meier A M 2013 *PhD Thesis* University of Colorado e-print: 1811.10040
- [16] Onorati E, Werner A H and Eisert J 2019 *Phys. Rev. Lett.* **123** 060501
- [17] Magesan E 2008 *PhD Thesis* University of Waterloo <https://uwaterloo.ca/handle/10012/6832>
- [18] Wallman J J 2017 Randomized benchmarking with gate-dependent noise *Quantum J.* **2** 47
- [19] Magesan E, Gambetta J M and Emerson J 2011 *Phys. Rev. Lett.* **106** 180504
- [20] Magesan E, Gambetta J M and Emerson J 2012 *Phys. Rev. A* **85** 042311
- [21] Wallman J J and Flammia S T 2014 *New J. Phys.* **16** 103032
- [22] Franca D S and Hashagen A L 2018 *J. Phys. A: Math. Theor.* **51** 215508
- [23] Proctor T J, Carignan-Dugas A, Rudinger K, Nielsen E, Blume-Kohout R and Young K 2018 arXiv:1807.07975
- [24] Merkel S T, Pritchett E J and Fong B H 2018 arXiv:1804.05951
- [25] Harrow A W and Low R A 2009 *Commun. Math. Phys.* **291** 257
- [26] Gottesman D 1998 Group 22: *Proc. of the 22nd Int. Coll. on Group Theoretic Methods in Physics* eprint: quant-ph/9807006
- [27] Helsen J, Xue X, Vandersypen L M K and Wehner S 2019 *Nature Phys. J.: Quantum Inf.* **5** 71
- [28] Onorati E, Buerschaper O, Kliesch M, Brown W, Werner A H and Eisert J 2017 *Commun. Math. Phys.* **355** 905
- [29] Brandao F G S L, Harrow A W and Horodecki M 2016 *Commun. Math. Phys.* **346** 397
- [30] Vermersch B, Elben A, Dalmonte M, Cirac J I and Zoller P 2018 *Phys. Rev. A* **97** 023604
- [31] Roberts D A and Yoshida B 2017 *J. High Energy Phys.* **JHEP04(2017)121**
- [32] Srednicki M 1994 *Phys. Rev. E* **50** 888
- [33] Brown W G, Santos L F, Starling D J and Viola L 2008 *Phys. Rev. E* **77** 021106
- [34] Lashkari N, Stanford D, Hastings M, Osborne T and Hayden P 2013 *J. High Energy Phys.* **JHEP04(2013)022**
- [35] Guhr T and Weidenmüller H A 1990 *Ann. Phys.* **199** 412
- [36] Swingle B, Bentsen G, Schleier-Smith M and Hayden P 2016 *Phys. Rev. A* **94** 040302
- [37] Hayden P and Preskill J 2007 *J. High Energy Phys.* **JHEP09(2007)120**
- [38] Marchildon L 2002 *Quantum Mechanics: From Basic Principles to Numerical Methods and Applications* Symmetry of the Hamiltonian (Amsterdam: Elsevier) pp 275–304
- [39] Abd El-Hady A, Abul-Magd A Y and Simbel M H 2002 *J. Phys. A: Math. Gen.* **35** 2361–72
- [40] Guhr T, Müller-Groaling A and Weidenmüller H A 1998 *Phys. Rep.* **299** 4
- [41] Blümel R and Smilansky U 1992 *Phys. Rev. Lett.* **69** 217
- [42] Kitaev A Y, Shen A H and Vedral M N 2002 *Classical and Quantum Computation* (Graduate Studies in Mathematics vol 47) (Providence, RI: American Mathematical Society)
- [43] Dankert C, Cleve R, Emerson J and Livine E 2009 *Phys. Rev. A* **80** 012304
- [44] Dai F and Xu Y 2013 *Cubature Formulas on Spheres Approximation Theory and Harmonic Analysis on Spheres and Balls* (New York, NY: Springer) 127–53
- [45] Harrow A W and Mehraban S 2018 arXiv:1809.06957
- [46] Gattner M, Bohnet J G, Safavi-Naini A, Wall M L, Bollinger J J and Rey A M 2017 *Nat. Phys.* **13** 781
- [47] Li J, Fan R, Wang H, Ye B, Zeng B, Zhai H, Peng X and Du J 2017 *Phys. Rev. X* **7** 031011
- [48] Lesovik G B, Sadovskyy I A, Suslov M V, Lebedev A V and Vinokur V M 2019 *Sci. Rep.* **9** 4396
- [49] Kim K *et al* 2011 *New J. Phys.* **13** 105003
- [50] Lanyon B P *et al* 2011 *Science* **334** 57

- [51] Britton J W, Sawyer B C, Keith A C, Wang C C J, Freericks J K, Uys H, Biercuk M J and Bollinger J J 2012 *Nature* **484** 489
- [52] Porras D and Cirac J I 2004 *Phys. Rev. Lett.* **92** 207901
- [53] Richerme P, Gong Z-X, Lee A, Senko C, Smith J, Foss-Feig M, Michalakis S, Gorshkov A V and Monroe C 2014 *Nature* **511** 198
- [54] Jurcevic P, Lanyon B P, Hauke P, Hempel C, Zoller P, Blatt R and Roos C F 2014 *Nature* **511** 202
- [55] Bernien H *et al* 2017 *Nature* **551** 579
- [56] The MathWorks, Inc. 2018 *Nonlinear Least Squares Fitting Tool, Matlab Curve Fitting Toolbox* (Natick, MA: The MathWorks)
- [57] Proctor T, Rudinger K, Young K, Sarovar M and Blume-Kohout R 2017 *Phys. Rev. Lett.* **119** 130502
- [58] Brown W and Fawzi O 2013 arXiv:1210.6644
- [59] Daley A J 2014 *Adv. Phys.* **63** 77
- [60] Hunter-Jones N 2019 arXiv:1905.12053
- [61] Low R A 2010 *PhD Thesis* University of Bristol e-print: 1006.5227
- [62] Nielsen M A 2002 *Phys. Rev. Lett. A* **303** 249
- [63] Issai S 1905 *Neue Begründung der Theorie der Gruppencharaktere* New foundation for the theory of group characters (Berlin: Sitzungsberichte der Königlich Preußischen Akademie der Wissenschaften zu) pp 406–32
- [64] Dankert C 2005 *PhD Thesis* University of Waterloo e-print: quant-ph/0512217
- [65] Gross D, Audenaert K and Eisert J 2007 *J. Math. Phys.* **48** 052104