

GATE-EFFICIENT DISCRETE SIMULATIONS OF CONTINUOUS-TIME QUANTUM QUERY ALGORITHMS

DOMINIC W. BERRY

*Department of Physics and Astronomy, Macquarie University
 Sydney, NSW 2109, Australia*

RICHARD CLEVE

*David R. Cheriton School of Computer Science and Institute for Quantum Computing, University of Waterloo
 Waterloo, Ontario N2L 3G1, Canada*

SEVAG GHARIBIAN

*Electrical Engineering & Computer Sciences, University of California
 Berkeley, California 94720, USA*

Received November 21, 2012

Revised May 15, 2013

We show how to efficiently simulate continuous-time quantum query algorithms that run in time T in a manner that preserves the query complexity (within a polylogarithmic factor) while also incurring a small overhead cost in the total number of gates between queries. By small overhead, we mean T within a factor that is polylogarithmic in terms of T and a cost measure that reflects the cost of computing the driving Hamiltonian. This permits *any* continuous-time quantum algorithm based on an efficiently computable driving Hamiltonian to be converted into a gate-efficient algorithm with similar running time.

Keywords: Quantum computation, quantum query complexity

Communicated by: R Jozsa & R Wolf

1 Introduction and Summary of Result

The standard quantum query model can be represented as an oracle that performs the unitary operation $|j, k\rangle \mapsto |j, k \oplus x_j\rangle$, where $x_1 x_2 \dots x_L \in \{0, 1\}^L$ is the data, and \oplus indicates modular addition. A convenient representation of the oracle is given by removing the ancilla, and having the oracle give a phase shift, so the unitary operation for the oracle, Q , acts as $Q|j\rangle = (-1)^{x_j}|j\rangle$. The fractional query model is a natural variant of this, where the operation is $Q^\lambda|j\rangle = (-1)^{\lambda x_j}|j\rangle = e^{i\pi\lambda x_j}|j\rangle$, and λ may be taken to be arbitrarily small (but positive). In the fractional query model, each size- λ query is taken to have cost λ . The fractional query model potentially provides more power than the standard query model, because additional unitary operations (which are independent of x_j) can be performed in between the fractional queries.

Informally, the continuous-time query model [1] arises from the fractional query model in the limit as λ approaches zero. More formally, in the continuous-time query model, the oracle operation is replaced with an oracle Hamiltonian, H_Q , which acts as $H_Q|j\rangle = x_j|j\rangle$. Evolving

under this Hamiltonian for time π would result in a full discrete query. The additional operations are replaced with a *driving Hamiltonian* independent of x_j , which we denote H (H may be time-dependent). The algorithm then becomes Hamiltonian evolution with the sum of the oracle and driving Hamiltonians, and the complexity is quantified by the time of evolution. The continuous-time and fractional query models are equivalent in the sense that each can simulate the other (to any desired level of accuracy) with the same query cost. For example, a continuous-time query algorithm can be approximated using fractional queries via a Lie-Trotter formula [2]. The continuous-time and discrete query models are also effectively equivalent, in that one can convert from one to the other with at most a polylogarithmic overhead in the query cost [2].

Presently, we are concerned not just with the query cost, but with the cost in terms of the number of additional gates and ancilla qubits needed. We show that any continuous-time quantum query algorithm whose total query time is T and whose driving Hamiltonian is implementable with G elementary gates (in a sense defined in Section 3) can be simulated by a discrete-query quantum algorithm using the following resources:

- $O(T \log T / \log \log T)$ queries
- $O(TG \log(T) + T \log^3(\|H\|T))$ elementary gates [or $O(TG \log(T) + TG^3)$ in terms of just T and G]
- $O(\log^3(\|H\|T))$ qubits of space [or $O(G^3)$].

This extends the previous result [2] where the query cost is the same, but where the orders of the second and third resource costs are at least $T^2 \text{polylog } T$ and $T \text{polylog } T$ respectively. The present result can also be compared with the result [3] where the query cost is superior to ours, $O(T)$ (which is asymptotically optimal), but whose methodology does not (as far as we know) yield an efficient gate construction from an efficiently implementable driving Hamiltonian.

Another advantage of our result is that it provides an exponential improvement in the scaling (of the number of gates and ancilla qubits) with $\|H\|$ over that in [2]. Here the number of gates is polylogarithmic in $\|H\|$, whereas it is superlinear in $\|H\|$ in [2]. This is important, as the norm of the driving Hamiltonian can potentially be large.

2 Significance to Quantum Computation

The continuous-time query model is an important tool for designing algorithms, and for example yielded the algorithm for AND-OR tree evaluation [4]. The difficulty with continuous-time quantum algorithms is that, in order to implement them on quantum computers, these abstract query algorithms need to be translated into concrete algorithms with subroutines substituted for the black-box queries.^a In these circumstances, what matters is the total gate complexity, which can be large if the cost of the operations performed between the queries is large, even if the number of queries is small. The contribution of our result is that it provides

^aA query is typically *not* something that could be physically implemented directly via continuous-time Hamiltonian evolution, as in an analog quantum computer. A query corresponds to the coherent evaluation of a classical function on several qubits, and requires several quantum gates to implement, regardless of whether it is a full query or a fractional query.

a *systematic* way to obtain a gate-efficient discrete-query algorithm from *any* continuous-time query algorithm where the driving Hamiltonian can be efficiently implemented. That is, whenever the implementation cost of the driving Hamiltonian is small, the total gate complexity is not much more than the query complexity times the cost of implementing each query.

Consider applying the continuous-time quantum algorithm in [4] for AND-OR tree evaluation to evaluate expressions of the form

$$\exists x_1 \forall x_2 \exists x_3 \cdots S x_L f(x_1, x_2, \dots, x_L), \quad (1)$$

where one is given a polynomial (in L) size circuit implementation of $f : \{0, 1\}^L \rightarrow \{0, 1\}$. The symbol S represents \forall for even L , and \exists for odd L . This corresponds to evaluating a balanced binary AND-OR tree of size $N = 2^L$. A continuous-time query algorithm achieving time $O(\sqrt{N})$ cannot be simulated directly from f , because a small λ -fractional query to f cannot be computed at cost proportional to λ ; the algorithm must be efficiently translated into the discrete-query framework to be implementable. But if we substitute the parameters into the simulation in [2], we obtain a gate cost of order $N \text{polylog } N$ (losing the square-root speedup) and consume order $\sqrt{N} \text{polylog } N$ qubits of space. The simulation in [3] does not appear to yield any bounds less than $O(N)$ on the gate cost. However, our present simulation results in $N^{1/2+o(1)}$ gates and $O(\text{polylog } N)$ space (using the fact that the driving Hamiltonian in [4] can be implemented with $N^{o(1)}$ gates). We remark that, for this particular example, a better simulation that is specific to AND-OR tree evaluation (that was discovered after [4]) is known [5, 6].

3 Precise Statement of Main Result

Prior to stating our main result, we give a precise definition of the *implementation cost* of a Hamiltonian acting on l qubits, which is the cost of realising the unitary operation corresponding to evolution under the Hamiltonian from a start time to a finish time. A preliminary idealised definition is as a unitary operation with the following properties. It acts on three registers: a *start time*, a *finish time* and an l -qubit *state*. For any start and finish times t_s and t_f , and any l -qubit state $|\psi\rangle$, the unitary operation maps $|t_s\rangle|t_f\rangle|\psi\rangle$ to $|t_s\rangle|t_f\rangle|\psi'\rangle$, where $|\psi'\rangle$ is the state that results when $|\psi\rangle$ evolves under H from time t_s to time t_f . Assuming that all three registers are finite-dimensional, this can be denoted as a gate as in Fig. 1. We will not require the unitary to be implemented perfectly. We introduce a precision parameter ε' , and permit the unitary evolution to be approximated within ε' . This leads to the following definition.

Definition 1 *Let H be a Hamiltonian acting on l qubits. Define H to be implementable within precision ε' with G gates if the following unitary operation can be implemented within precision ε' with G elementary gates. These elementary gates can be taken to be any unitary gates acting on at most two qubits. The unitary acts on three registers: a start time and finish time, and l qubits set to the initial state. The unitary maps $|t_s\rangle|t_f\rangle|\psi\rangle$ to $|t_s\rangle|t_f\rangle|\psi'\rangle$, where $|\psi'\rangle$ is the state that results when $|\psi\rangle$ evolves under H from time t_s to time t_f . By approximating within ε' , we mean with respect to the completely bounded norm.*

We are now ready to state our main result.

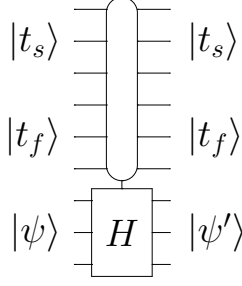


Fig. 1. Controlled evolution under Hamiltonian H , with start time t_s , finish time t_f , and target state $|\psi\rangle$.

Theorem 1 (Main) *Let $H(t)$ be a driving Hamiltonian that is approximately implementable within precision $O(1/T)$ using G gates. Then the continuous-time query algorithm can be simulated with constant error by a discrete-query quantum algorithm using the following resources:*

- $O(T \log T / \log \log T)$ queries
- $O(TG \log(T) + T \log^3(\|H\|T))$ elementary gates
- $O(\log^3(\|H\|T))$ qubits of space.

In particular, when G is $\text{polylog}(T)$, this is $\tilde{O}(T)$ queries, $\tilde{O}(T)$ elementary gates, and $\text{polylog}(T)$ qubits of space. The norm $\|H\|$ is taken to be $\|H\| := \sup_{t \in [0, T]} \|H(t)\|$ for time-dependent $H(t)$. Because the gate complexity scales linearly in G , we require the driving Hamiltonian to be simulatable efficiently in order for the simulation to be gate-efficient. If, for example, G scaled linearly in $\|H\|$, then the gate complexity would be linear in $\|H\|T$, which is similar to the complexity obtained by product formulae [7]. On the other hand, we have a lower bound of $G = \Omega(\log(\|H\|T))$ (see Section 4.6). As a result, we could express the gate complexity as $O(TG \log(T) + TG^3)$, and the number of qubits of space as $O(G^3)$.

The remaining sections explain our algorithm, with the proof of Theorem 1 in Section 6.

4 Compressed CGMSY Construction

We will summarise the construction in [2], and then show how to make it more efficient by compressing the control registers. Before doing so, we state the notation used throughout this paper.

Notation. We denote the set of linear operators acting on complex Euclidean space \mathcal{X} as $\mathcal{L}(\mathcal{X})$. The spectral norm of operator A is $\|A\| := \max\{\|A|v\rangle\|_2 : \|v\rangle\|_2 = 1\}$. The norm of time dependent operator $A(t)$ is given by $\|A\| = \sup_t \|A(t)\|$. The completely bounded norm, or diamond norm, of superoperator $\Phi : \mathcal{L}(\mathcal{X}) \mapsto \mathcal{L}(\mathcal{Y})$ is defined as $\|\Phi\|_\diamond = \|\Phi \otimes I_{\mathcal{L}(\mathcal{X})}\|_1$, where the superoperator trace-norm is given by $\|\Phi\|_1 = \max\{\|\Phi(X)\|_{\text{tr}} : X \in \mathcal{L}(\mathcal{X}), \|X\|_{\text{tr}} \leq 1\}$. All logarithms are taken to base 2. We define $[m] := \{1, \dots, m\}$. The tensor product of many zero computational basis states will be represented in compact form as $|0^\ell\rangle := |0\rangle^{\otimes \ell}$.

4.1 Overview of the CGMSY Construction [2]

Our result is obtained by simulating the construction in [2], but by representing some of the qubits in a highly compressed form. This compressed form was known by the authors of [2], but it was not known that all of the steps of the construction can be carried out within the compressed form—especially the *measurement of control qubits*.

The construction in [2] begins with a continuous-time query algorithm with total query cost T . The overall Hamiltonian for the continuous-time query algorithm is a sum of the oracle Hamiltonian and the driving Hamiltonian, so the evolution can be approximated via a Lie-Trotter decomposition. As above, it is assumed that the driving Hamiltonian can be simulated, and the evolution under the oracle Hamiltonian for a short time becomes a fractional-time query.

The total time T is partitioned into segments corresponding to time intervals of the form $[t_0, t_0 + 1/4]$, and with m of the Lie-Trotter time intervals within each segment. We call each length $1/4$ time interval a *segment*, to distinguish them from other time intervals considered. In each of the Lie-Trotter time intervals there is a fractional query of size $1/4m$. Here, m can be chosen as a power of two without loss of generality; we henceforth assume this is the case. In this work we consider the simulation of each of these segments.

Within each segment, there are m fractional queries which we wish to simulate. The method in [2] is to then, for each fractional query, use a *control* qubit that is in the state $\alpha|0\rangle + i\beta|1\rangle$. The unitary operation for the discrete oracle, Q , is then implemented, controlled by the control qubit. Given that the *target* system is initially in state $|\zeta\rangle$, the state after this controlled operation is

$$\alpha|0\rangle \otimes |\zeta\rangle + i\beta|1\rangle \otimes Q|\zeta\rangle. \quad (2)$$

Finally, a projection measurement with outcome $\alpha|0\rangle + \beta|1\rangle$ yields the state in the target system (omitting normalisation)

$$\alpha^2|\zeta\rangle + i\beta^2Q|\zeta\rangle. \quad (3)$$

The query Hamiltonian, H_Q , has values on the diagonal equal to x_j , whereas the discrete query unitary Q has values on the diagonal of $(-1)^{x_j} = 1 - 2x_j$. Therefore the Hamiltonian and unitary are related by $H_Q = (I - Q)/2$. The I only gives a global phase factor and can be ignored. Because Q is self-inverse, one obtains (omitting the phase factor)

$$e^{-iH_Q t} = \cos(t/2)I + i\sin(t/2)Q. \quad (4)$$

With $t = 1/4m$, one therefore obtains the correct operation via the above procedure if $\beta \approx 1/\sqrt{8m}$.

The number of calls to the oracle can then be reduced by, instead of considering controlled operations at each time step individually, considering them jointly within a segment. That is, considering the state of all control qubits together, for a given basis state the position of each 1 gives a time that Q is applied. As the only basis states with significant weighting are those with a small number of ones, we can allow a maximum number $k' \in O(\log(T)/\log \log(T))$ of applications of the oracle, with evolution under the driving Hamiltonian between them. That is, the positions of the ones in the control qubits control the time of evolution under the driving Hamiltonian.

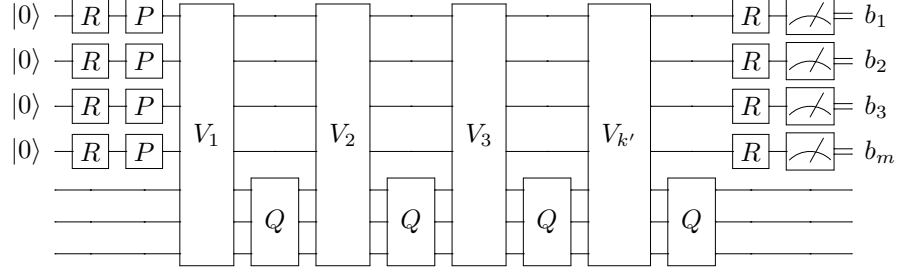


Fig. 2. The construction from Ref. [2] to simulate a segment corresponding to a time interval of length $1/4$.

This procedure from [2] is represented in Fig. 2. The operations P and R are designed to prepare the initial qubits and enable the final measurement, and are given by

$$P = \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}, \quad R = \begin{pmatrix} \alpha & \beta \\ \beta & -\alpha \end{pmatrix} \quad \text{with } \beta \approx 1/\sqrt{8m}. \quad (5)$$

The sequence of operations PR acting on $|0\rangle$ prepares $\alpha|0\rangle + i\beta|1\rangle$, and R followed by a computational basis measurement of $b_j = 0$ corresponds to the desired measurement outcome $\alpha|0\rangle + \beta|1\rangle$. The gates $V_1, \dots, V_{k'}$ are the unitaries corresponding to evolving the driving Hamiltonian for various time intervals specified by the control qubits: V_1 for the time interval from t_0 to the position of the first one in the control qubits; V_2 for the time interval delineated by the positions of the first and second ones in the control qubits; and so on. The simulation is successful if $b_1 = \dots = b_m = 0$. The probability of obtaining each $b_j = 0$ is at least $1 - 1/4m$, and there are m measurements, so the probability of successful simulation is at least $3/4$. The value $\beta^2 \approx 1/8m$ corresponds to a time interval of $1/4$. This time interval is chosen to ensure that the success probability is at least $3/4$.

In the case that the simulation is not successful, there are errors at times corresponding to the b_j that are equal to 1. Reference [2] shows how to correct unsuccessful instances. Since the errors are unitary operations, it is possible to undo the step that has just been performed, then redo it. To undo the step, one inverts the construction given in Fig. 2, but with each of the errors inverted. This inversion will also succeed with probability at least $3/4$. If this inversion does not succeed, then one attempts to undo it and then redo it, and so forth. This procedure corresponds to a biased random walk, where a step to the right (corresponding to a success) occurs with probability at least $3/4$, and a step to the left (corresponding to a failure) occurs with probability at most $1/4$. Overall success for this random walk is obtained when it advances one step to the right of its initial position.

That analysis continues to hold here without modification. The only subtlety is that we also need to account for the number of gates needed to perform the gates V_j . Each gate may need to be divided into a number of parts corresponding to the number of errors (ones) found. It is shown in Ref. [2] that the average number of ones is $O(1)$, so the *average* number of oracle queries is at most multiplied by a constant factor. Moreover, if the total number of oracle queries permitted is bounded by $O(1/\varepsilon_{\text{tot}})$ times the average value, then by the Markov bound, the probability is at least $1 - O(\varepsilon_{\text{tot}})$ that the overall correction procedure terminates within this bound [2]. Failure to terminate within the bound can be included in the ε_{tot}

allowable error. For the main result in Theorem 1, constant error is considered, so this does not alter the result.

When analysing the complexity due to correcting unsuccessful instances, another factor that needs to be considered is the additional complexity due to correcting the individual errors. The average of this complexity was denoted C_0 in Ref. [2], but an upper bound was not considered. As before, an upper bound equal to $O(1/\varepsilon_{\text{tot}})$ times the average value will not be exceeded with probability $1 - O(\varepsilon_{\text{tot}})$. Again this does not affect the result in Theorem 1 as constant error is considered. As a result of these considerations, when taking into account the corrections, the number of oracle queries and the number of additional elementary gates are at most multiplied by a constant factor. This means that the correction operations do not alter the scaling, and we do not need to consider them further.

The feature of the analysis in [2] that is most crucial for this work is that the state of the control registers $R^{\otimes m}|0^m\rangle = (\alpha|0\rangle + \beta|1\rangle)^{\otimes m}$ is highly “compressible” in that most of its amplitude is concentrated on basis states with low Hamming weight. A natural succinct representation of this state is in terms of the positions of the ones in binary. We first define such a succinct form precisely (Section 4.2). We then show how the above circuit can be simulated with the control qubits in their succinct form in these three stages: the initial stage (Sections 4.3 and 4.4), which is the construction of the state $R^{\otimes m}|0^m\rangle$; the intermediate stage (Section 4.5), where $P^{\otimes m}$ is applied to the control qubits and then the queries and driving operations occur; and the final stage (Section 5), which is where the control qubits are measured with respect to the basis $\{R^{\otimes m}|x\rangle : x \in \{0, 1\}^m\}$.

4.2 Succinct Representation of Control Qubits

We now propose a succinct encoding scheme which accurately reproduces low Hamming weight basis states. Specifically, consider the set of all m -bit strings whose Hamming weight is at most $k + 1$, where k is much smaller than m . The size of this set is bounded above by $(m + 1)^{k+1}$. Our encoding scheme utilizes a set of size $(m + 1)^{k+1}$ strings to accurately represent this space as follows. We use the notation $|x|$ to denote the Hamming weight of $x \in \{0, 1\}^*$. The value of k is chosen to ensure that the error due to omitting high Hamming weight components is no more than ε , and therefore can be taken as

$$k = \Theta\left(\frac{\log(1/\varepsilon)}{\log \log(1/\varepsilon)}\right). \quad (6)$$

We also use a slightly smaller value k' to ensure that the error is no more than ε' ; the relation between these primed variables is identical. The Hamming weight cutoff k is used to limit errors that occur repeatedly in the compressed measurement protocol. In contrast, Hamming weight cutoff k' is used to limit errors that only occur once. In particular, we limit the total number of controlled oracle calls to k' , because the error due to limiting the Hamming weight there only occurs once. We also limit the number of ones that are measured to k' . The Hamming weight cutoff k is used in our compressed encoding, as the error due to this cutoff will contribute multiple times.

Definition 2 Define the encoding scheme C_n^k on $|x\rangle$ for $x \in \{0, 1\}^n$, $|x| \leq k$ as follows. For

$x = 0^{s_1}10^{s_2}10^{s_3} \dots 0^{s_h}10^t$, where $h := |x|$, $h \leq k+1$ and $t = n - s_1 - \dots - s_h - h$,

$$C_n^k|x\rangle = |s_1, s_2, \dots, s_h, \underbrace{n, \dots, n}_{k+1-h}\rangle, \quad (7)$$

where $C_n^k|x\rangle \in (\mathbb{C}^{n+1})^{\otimes k+1}$. For $h > k+1$, $C_n^k|x\rangle$ encodes the positions of the first $k+1$ ones.

We have allowed the encoding to act upon n qubits for generality. We will initially use $n = m$ for state preparation, but will use $n < m$ when we break apart the encoding for use in the measurement. Note that the rotations R always use m , rather than n .

4.3 Initialization of Control Qubits in Alternative Encoding

We now show how to simulate the preparation of the state after operation $R^{\otimes m}$ (but before $P^{\otimes m}$, which is deferred to Section 4.5) in succinct form using the encoding of Definition 2. In order to achieve this, we efficiently prepare a state according to the following Theorem.

Theorem 2 For $n \leq m$ and k as in Eq. (6), it is possible to prepare an approximation, within trace distance $O(\varepsilon)$, of the state

$$|\Xi_n^k\rangle := \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} C_n^k|x\rangle + \mu|\nu'\rangle, \quad (8)$$

where $|\nu'\rangle \in (\mathbb{C}^{n+1})^{\otimes k+1}$ is orthogonal to all the kets arising in the sum, via a quantum circuit with complexity

$$O(k[\log m + \log \log(1/\varepsilon)]). \quad (9)$$

Before we prove this Theorem, there are a number of intermediate results that we need to prove. The most important property of the state $|\Xi_n^k\rangle$ is that the inner products are the same as for the uncompressed state.

Lemma 1 For $n \leq m$ and $|x| \leq k$,

$$\langle x|(C_n^k)^\dagger|\Xi_n^k\rangle = \langle x|R^{\otimes n}|0\rangle. \quad (10)$$

For k as in Eq. (6), $\mu^2 \in O(\varepsilon)$, and the norm of the component of $R^{\otimes n}|0\rangle$ with Hamming weight greater than k is $O(\varepsilon)$.

Proof. There are n control qubits, each of which is rotated by R , in order to give the state $(\alpha|0\rangle + \beta|1\rangle)^{\otimes n}$, where $\beta \in \Theta(1/\sqrt{m})$. The amplitudes of terms in this superposition decrease factorially with Hamming weight, and in particular, one can write

$$\begin{aligned} R^{\otimes n}|0\rangle &= (\alpha|0\rangle + \beta|1\rangle)^{\otimes n} \\ &= \sum_{x \in \{0,1\}^n} \alpha^{n-|x|} \beta^{|x|} |x\rangle \\ &= \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} |x\rangle + \sum_{\substack{x \in \{0,1\}^n \\ |x| > k}} \alpha^{n-|x|} \beta^{|x|} |x\rangle \\ &= \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} |x\rangle + \mu|\nu\rangle. \end{aligned} \quad (11)$$

One then obtains that

$$\langle x | R^{\otimes n} | 0 \rangle = \alpha^{n-|x|} \beta^{|x|} = \langle x | (C_n^k)^\dagger | \Xi_n^k \rangle, \quad (12)$$

as required. On the last line of Eq. (11), $|\nu\rangle$ is orthogonal to every basis state in the sum that precedes it, and therefore μ here is the same as in Eq. (8). Moreover $\mu^2 \in 1/2^{O(k)} k!$, so using k as in Eq. (6), $\mu^2 \in O(\varepsilon)$. This means that the norm of the component with Hamming weight greater than k is $O(\varepsilon)$ \square .

We now show how to construct an approximation (within distance ε) of the state in Eq. (8) using $\text{poly}(k, \log m)$ gates. Note that, to accomplish this, we must avoid any approach based on first constructing the expanded state in Eq. (11) then applying C_n^k , since this would immediately entail order m gates. Our efficient approach is to first prepare a state similar to Eq. (11) using a slightly different encoding scheme than C_n^k , denoted $B_{n,q}^k$. We then postprocess the state so that the encoding is changed from $B_{n,q}^k$ to C_n^k [i.e. Eq. (8)].

We now introduce the encoding $B_{n,q}^k$ by explicit construction. Specifically, it is based on the *exponential superposition* state, which can be efficiently constructed.

Lemma 2 *The exponential superposition state*

$$|\phi_q\rangle := \sum_{s=0}^{q-1} \beta \alpha^s |s\rangle + \alpha^q |q\rangle, \quad (13)$$

where $q = 2^r$, can be prepared using $O(r)$ elementary operations.

Proof. The state $|\phi_q\rangle$ is very simple to prepare as follows. Define the unitary matrix

$$M(\gamma) := \frac{1}{\sqrt{1+\gamma^2}} \begin{pmatrix} 1 & -\gamma \\ \gamma & 1 \end{pmatrix}. \quad (14)$$

Note that

$$\begin{aligned} & M(\alpha^{2^{r-1}}) \otimes \cdots \otimes M(\alpha^2) \otimes M(\alpha) |0^r\rangle \\ &= \frac{\beta}{\sqrt{1-\alpha^{2q}}} (|0 \dots 00\rangle + \alpha |0 \dots 01\rangle + \alpha^2 |0 \dots 10\rangle + \cdots + \alpha^{q-1} |1 \dots 11\rangle) \\ &= \frac{1}{\sqrt{1-\alpha^{2q}}} \sum_{s=0}^{q-1} \beta \alpha^s |s\rangle. \end{aligned} \quad (15)$$

Therefore, a circuit that maps $|0^{r+1}\rangle$ to $|\phi_q\rangle$ can be obtained by first applying a one-qubit gate on the first qubit to put it in state $\sqrt{1-\alpha^{2q}}|0\rangle + \alpha^q|1\rangle$, and then applying a sequence of controlled- $M(\alpha^{2^j})$ gates (each controlled by the first qubit being in state $|0\rangle$) to create the state

$$\begin{aligned} & \beta |0\rangle (|0 \dots 00\rangle + \alpha |0 \dots 01\rangle + \alpha^2 |0 \dots 10\rangle + \cdots + \alpha^{q-1} |1 \dots 11\rangle) + \alpha^q |1\rangle |0 \dots 00\rangle \\ &= \beta (|00 \dots 00\rangle + \alpha |00 \dots 01\rangle + \alpha^2 |00 \dots 10\rangle + \cdots + \alpha^{q-1} |01 \dots 11\rangle) + \alpha^q |10 \dots 00\rangle \\ &= |\phi_q\rangle. \end{aligned} \quad (16)$$

As there is one single-qubit gate, and r two-qubit controlled gates used, the total number of gates is $O(r)$ \square .

The reason why state $|\phi_q\rangle$ is useful is because, for $q \geq n$, $|\phi_q\rangle^{\otimes k+1}$ yields a state similar to Eq. (8). The encoding, which we will call $B_{n,q}^k|x\rangle$, is slightly different than $C_n^k|x\rangle$, but can be efficiently translated into $C_n^k|x\rangle$ with some “clean-up” operations. Specifically, the encoding is as defined below.

Definition 3 Define the encoding scheme $B_{n,q}^k$ on $|x\rangle$ for $x \in \{0,1\}^n$, $|x| \leq k$ as follows. For $x = 0^{s_1}10^{s_2}10^{s_3} \dots 0^{s_h}10^t$, where $h := |x|$, $h \leq k$ and $t = n - s_1 - \dots - s_h - h$,

$$B_{n,q}^k|x\rangle = |s_1, \dots, s_h\rangle \left(\sum_{j=0}^{q-t-1} \alpha^j \beta |j+t\rangle + \alpha^{q-t} |q\rangle \right) |\phi_q\rangle^{\otimes k-h}, \quad (17)$$

where $B_{n,q}^k|x\rangle \in (\mathbb{C}^{q+1})^{\otimes k+1}$.

The state is then given as in the following theorem.

Theorem 3 For $q \geq n$, the states $B_{n,q}^k|x\rangle$ for $|x| \leq k$ are orthonormal and

$$|\phi_q\rangle^{\otimes k+1} = \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} B_{n,q}^k|x\rangle + \mu|\nu'\rangle, \quad (18)$$

for some $|\nu'\rangle$ orthogonal to all $B_{n,q}^k|x\rangle$ for $|x| \leq k$.

Proof. The state $|\phi_q\rangle^{\otimes k+1}$ is a superposition of (computational) basis states of the form $|s_1, \dots, s_{k+1}\rangle$, where $s_1, \dots, s_{k+1} \in \{0, 1, \dots, q\}$. Intuitively, it is useful to think of each such basis state as an encoding of a binary string $0^{s_1}10^{s_2}1 \dots 0^{s_{k+1}}1$ (whose Hamming weight is $k+1$ and length is $s_1 + \dots + s_{k+1} + k+1$). We will show that these basis states can be naturally partitioned into equivalence classes: one for each prefix $x \in \{0,1\}^n$ with $|x| \leq k$, and one for all the remaining basis states.

Let $x \in \{0,1\}^n$ with $h = |x| \leq k$ be of the form $x = 0^{s_1}10^{s_2}10^{s_3} \dots 0^{s_h}10^t$. Consider the set P_x that consists of all $|s'_1, s'_2, \dots, s'_{k+1}\rangle$ that are encodings of strings whose n -bit prefix is x . The set P_x consists of all $|s'_1, s'_2, \dots, s'_{k+1}\rangle$ such that $(s'_1, s'_2, \dots, s'_h) = (s_1, s_2, \dots, s_h)$, $s'_{h+1} \in \{t, \dots, q\}$, and $s'_{h+2}, \dots, s'_{k+1} \in \{0, \dots, q\}$. It follows that the sum of all the terms in the superposition

$$|\phi_q\rangle^{\otimes k+1} = \sum_{s'_1=0}^q \sum_{s'_2=0}^q \dots \sum_{s'_{k+1}=0}^q \alpha^{s'_1+s'_2+\dots+s'_{k+1}} \beta^{|\{\ell | s'_\ell < q\}|} |s'_1, s'_2, \dots, s'_{k+1}\rangle \quad (19)$$

that correspond to elements of P_x is

$$\begin{aligned} & \alpha^{s_1} \beta \dots \alpha^{s_h} \beta |s_1, \dots, s_h\rangle \left(\sum_{j=t}^{q-1} \alpha^j \beta |j\rangle + \alpha^q |q\rangle \right) |\phi_q\rangle^{\otimes k-h} \\ &= \alpha^{s_1} \beta \dots \alpha^{s_h} \beta \alpha^t |s_1, \dots, s_h\rangle \left(\sum_{j=0}^{q-t-1} \alpha^j \beta |j+t\rangle + \alpha^{q-t} |q\rangle \right) |\phi_q\rangle^{\otimes k-h} \\ &= \alpha^{n-|x|} \beta^{|x|} B_{n,q}^k|x\rangle, \end{aligned} \quad (20)$$

which is the appropriate weighting for $B_{n,q}^k|x\rangle$ in the sum in Eq. (18).

Thus, the basis states in the superposition in Eq. (19) corresponding to encodings of strings $x \in \{0,1\}^n$ of Hamming weight at most k can be grouped into equivalence classes P_x .

What about the remaining terms in $|\phi_q\rangle^{\otimes k+1}$ which do not fall in any P_x ? These are the $|s_1, \dots, s_{k+1}\rangle$ where $s_1 + \dots + s_{k+1} + k + 1 \leq n$. Therefore, we can set

$$\mu|\nu'\rangle = \sum_{s_1 + \dots + s_{k+1} + k + 1 \leq n} \alpha^{s_1 + \dots + s_{k+1}} \beta^{k+1} |s_1, \dots, s_{k+1}\rangle, \quad (21)$$

where $\mu \in \mathbb{R}$ is chosen so that $|\nu'\rangle$ is normalised. All the $B_{n,q}^k|x\rangle$ and $|\nu'\rangle$ are mutually orthogonal since they are constructed from a partition of the basis states \square .

4.4 Converting from the B Encoding to the C Encoding

We have thus far shown how to prepare states in the encoding $B_{n,q}^k$. As mentioned above, we can now convert from the encoding $B_{n,q}^k$ to our desired encoding C_n^k . We find that it is possible to efficiently convert between the encodings with the efficiency given by the following lemma.

Lemma 3 *For $n \leq m$ it is possible to convert from the state $|\phi_q\rangle^{\otimes k+1}$ to $|\Xi_n^k\rangle$ within error $O(\varepsilon)$ for $\log q \in \Theta(\log m + \log \log(1/\varepsilon))$ and q a power of two, and using*

$$O(k[\log m + \log \log(1/\varepsilon)]) \quad (22)$$

elementary operations.

Proof. This is achieved by “cleaning up” the registers that follow register $h = |x|$ in $B_{n,q}^k|x\rangle$ [compare Eq. (7) with Eq. (17)]. The difference is that, instead of these registers being in the state $|n\rangle$, they are in the state $|\phi_q\rangle$ (for registers $h + 2$ to $k + 1$). Register $h + 1$ is in a state that is similar to $|\phi_{q-t}\rangle$, except that the basis states are shifted by t . Therefore, we need a way of converting these registers to the state $|n\rangle$. However, this conversion depends on both h and t , so we first need these quantities. We will first give a simplified explanation, then expand on the technical details. To determine h and t , we compute the prefix sums

$$|s_1\rangle|s_2\rangle \cdots |s_{k+1}\rangle \mapsto |s_1 + 1\rangle|s_1 + s_2 + 2\rangle \cdots |s_1 + s_2 + \dots + s_{k+1} + k + 1\rangle. \quad (23)$$

This gives the *absolute* positions of the ones. The value of h can be determined by finding the first register with a value larger than n (which would give a position for a one past the end of the string).

Now we can identify register $h + 1$. For this register, we wish to subtract t , so that the state of this register [as in Eq. (17)] becomes $|\phi_{q-t}\rangle$. At this stage we have computed the prefix sums, and subtracting $n + 1$ from this modified register gives the same result as subtracting t from the unmodified register. That is, we do not need to explicitly compute t to subtract it, because it is obtained implicitly in the prefix sum. For all the other registers we then undo the prefix sums.

At this stage we have h in an ancilla, and we have subtracted t from register $h + 1$. Now we can undo the procedure to prepare $|\phi_q\rangle$ in registers $h + 1$ to $k + 1$. Register $h + 1$ is actually in state $|\phi_{q-t}\rangle$ rather than $|\phi_q\rangle$, but it is a good approximation of state $|\phi_q\rangle$. Therefore the inverse preparation yields states $|0\rangle$ in registers $h + 1$ to $k + 1$, with this being approximate for register $h + 1$. It is trivial to convert $|0\rangle$ to $|n\rangle$, then uncompute the value of h in the ancilla register. This then completes the conversion of the encoding.

In summary the overall procedure is as follows.

1. Compute the prefix sums.
2. Compute $h = |x|$ in an ancilla register.
3. Uncompute the prefix sums for registers other than $h + 1$, and subtract $n + 1$ from register $h + 1$.
4. Invert the procedure to prepare $|\phi_q\rangle$ from $|0\rangle$ on registers $h + 1$ to $k + 1$, and swap register $h + 1$ with the error flag register.
5. Flip one qubit on registers $h + 1$ to $k + 1$ to change $|0\rangle$ to $|n\rangle$.
6. Uncompute h in the ancilla register.

Next we explain the technical details, including the error flag register. When computing the prefix sums, we can first consider the case of low-Hamming weight strings with $h \leq k$. For the first h registers the result is at most n , whereas for register $h + 1$, the result is (coherently) more than n . To prevent the value in register $h + 1$ wrapping around modulo n , we instead expand the registers to dimension $n + q + 2$, and perform the computations modulo $n + q + 2$. Because the value in register $h + 1$ is no more than that in h (which is $\leq n$) plus $q + 1$, the value is $\leq n + q + 1$, and does not wrap around modulo $n + q + 2$. The values in registers $h + 2$ to $k + 1$ may wrap around, but this does not affect the calculation. This covers steps 1 and 2 above.

Next, considering step 3, the value in register $h + 1$ will be

$$s_1 + \dots + s_h + s_{h+1} + h + 1 = n - t + s_{h+1} + 1. \quad (24)$$

We aim to obtain $s_{h+1} - t$ in this register. If we had computed the value of t , we could uncompute the prefix sums, then subtract t . However, it is obvious from Eq. (24) that we can just subtract $n + 1$ instead. Note that this is the first register that is larger than n , so subtracting $n + 1$ does not result in a negative number. We also need to uncompute the prefix sums for all registers other than register $h + 1$. This can be achieved by working backwards from register $k + 1$ to $h + 2$ uncomputing prefix sums, subtracting $n + 1$ from register $h + 1$, then uncomputing prefix sums from register h back to 1.

Next we consider the inverse preparation in step 4. At this stage, we have subtracted t from register $h + 1$ yielding the exponential state

$$|\phi_{q-t}\rangle = \sum_{s=0}^{q-t-1} \beta \alpha^s |s\rangle + \alpha^{q-t} |q-t\rangle. \quad (25)$$

By choosing q to be sufficiently large, $|\phi_{q-t}\rangle$ is close to $|\phi_q\rangle$, and inverting the procedure for preparing $|\phi_q\rangle$ yields an accurate approximation of $|0^{r+1}\rangle$. To be more precise, note that $\langle \phi_{q-t} | \phi_q \rangle = 1 - (1 - \beta) \alpha^{2(q-t)}$. Therefore, we have $\langle \phi_{q-t} | \phi_q \rangle \geq 1 - \varepsilon$ if $q \geq m + (1/\beta^2) \log(1/\varepsilon)$. To achieve this, $|\phi_q\rangle$ need only consist of $\log(m + 1/\beta^2) + \log \log(1/\varepsilon) + O(1)$ qubits. In particular, in our context where $\beta = \Theta(1/\sqrt{m})$, the number of qubits is $\log m + \log \log(1/\varepsilon) + O(1)$, so the precision scales double exponentially with the number of additional qubits beyond $\log m$.

This approximate step could alternatively be performed using the state preparation procedure of Grover and Rudolph [8]. Another alternative is to use amplitude amplification to ensure that the register is set to zero correctly. These alternatives would also not be exact, because they would require the coherent calculation of trigonometric functions.

It is convenient for the analysis to swap register $h + 1$ with an “error flag” register that has been prepared in the $|0\rangle$ state. Then, if this register is measured as not zero, it flags that the clean-up operation has not occurred properly. On the other hand, register $h + 1$ is exactly $|0\rangle$.

We also need to take account of the action of the conversion procedure on the state $|\nu'\rangle$. This state is a superposition of basis states $|s_1, \dots, s_{k+1}\rangle$, where $s_1 + \dots + s_{k+1} + k + 1 \leq n$. This means that, when we compute the prefix sums, the last register will *not* be $> n$. In this case, we can set $h = k + 1$, and then make no changes to the other registers in steps 3 to 5 for this value of h . This means that $|\nu'\rangle$ is unchanged. The exact form of this state is unimportant, because it corresponds to an error. However, $|\nu'\rangle$ is a superposition of strings of Hamming weight $h + 1$ encoded using C_n^k , and remains so under the conversion.

In summary, the overall preparation procedure is to prepare the state $|\phi_q\rangle^{\otimes k+1}$, then perform the clean-up procedure consisting of steps 1 to 6 above. By choosing $\log q \in \Theta(\log m + \log \log(1/\varepsilon))$ (for q a power of two), this then yields the state (8) within distance $O(\varepsilon)$. Our circuit has size given by Eq. (22). The final state has no values in its registers larger than n , so it can be stored in registers of dimension $n + 1$, though higher dimensions are required in intermediate steps \square .

To prepare the state, we have started with all qubits of registers in the state $|0\rangle$. It is convenient to start with these registers in the state $|n\rangle$, flip one qubit in each register to give $|0\rangle$, then perform the preparation procedure as described above. Then we are mapping the state $C_n^k|0^n\rangle$ (which is the state $|n\rangle^{\otimes k+1}$) to the succinct representation of $(\alpha|0\rangle + \beta|1\rangle)^{\otimes n}$ as defined in Eq. (8).

Finally we can easily prove Theorem 2 using the above results.

Proof. (of Theorem 2) First, using Lemma 2 we can prepare the state $|\phi_q\rangle$ using $O(r)$ elementary operations, and can therefore prepare $|\phi_q\rangle^{\otimes k+1}$, for $\log q \in \Theta(\log m + \log \log(1/\varepsilon))$ using

$$O(k[\log m + \log \log(1/\varepsilon)]) \tag{26}$$

elementary operations. Then, by Lemma 3 we can convert to the state $|\Xi_n^k\rangle$ within $O(\varepsilon)$ with the same complexity \square .

4.5 Phase Gates, Queries and Driving Operations

Applying the phase gates, $P^{\otimes m}$, to the control qubits in their succinct representation is straightforward because $P^{\otimes m}|x\rangle = i^{|x|}|x\rangle$. We need only compute $|x|$ in an ancilla register, apply $|s\rangle \mapsto i^s|s\rangle$, and then uncompute $|x|$ in the ancilla.

To apply the driving operations, we note that our definition of driving Hamiltonian implementation fits perfectly in this context, once we compute the prefix sums to give the positions of the ones, as in Eq. (23). In the compressed representation, V_1 is the implementation of the driving Hamiltonian with t_s hardwired to 0 and t_f controlled by the first register. V_2 is the implementation with t_s controlled by the first register and t_f controlled by the second register, and so on. At the end, the prefix sums can be uncomputed.

4.6 The Value of m Needed

In the CGMSY construction the number of fractional queries m comes from breaking up the evolution under the oracle and the driving Hamiltonian via a product formula. To obtain error bound by ε_{tot} with evolution over time T and driving Hamiltonian with norm $\|H\|$, the number of time intervals needed in a Lie-Trotter-Suzuki product formula for constant Hamiltonian H is $O(\|H\|T(\|H\|T/\varepsilon_{\text{tot}})^\delta)$ [7]. For the CGMSY construction the intervals need to be of equal size, which restricts δ to $1/2$.

For time-dependent Hamiltonians, the complexity of Lie-Trotter-Suzuki product formulae will depend on the magnitude of the derivatives of H when one is sampling the Hamiltonians at different times [9]. The situation we have here is somewhat different, because we assume that the evolution under the time-dependent driving Hamiltonian can be implemented. In this case, the error does not depend on the time derivative, and the error for a short time interval δt can be bounded as $\|H\|\delta t^2$ (this is easily derived from Eq. (2.3) of Ref. [10]). Hence the number of intervals to limit the overall error to $O(\varepsilon_{\text{tot}})$ need be no greater than $O(\|H\|T^2/\varepsilon_{\text{tot}})$. The number of intervals in one CGMSY segment of length $O(1)$ is therefore $m = O(\|H\|T/\varepsilon_{\text{tot}})$.

Another question is the precision that the time needs to be specified to in order to limit the overall error to ε_{tot} . It is easily shown that the error in the time needs to be bound as $O(\varepsilon'/\|H\|)$ in order to limit the error in a single operation to ε' . If the time is being specified on the interval $[0, T]$, then the number of bits needed for the time is $\lceil \log(\|H\|T/\varepsilon') \rceil$. Because there are $O(1)$ controlled Hamiltonian evolutions in each CGMSY segment, we need $\varepsilon' = O(\varepsilon_{\text{tot}}/T)$. This gives the number of bits for the time as $\log(\|H\|T^2/\varepsilon_{\text{tot}}) + O(1)$ (where the constant $O(1)$ is because ε' may have a constant of proportionality with $\varepsilon_{\text{tot}}/T$).

This result is consistent with the value of m used, because $\log(\|H\|T^2/\varepsilon_{\text{tot}}) + O(1)$ bits are needed to specify an integer from 0 to $O(mT)$. In the CGMSY construction, a superposition over the m time intervals is used, so the number of qubits needed is $\lceil \log m \rceil$. The number of the CGMSY segment also needs to be stored, but that can be stored in $O(\log T)$ classical bits.

One can use the number of bits for the time to place a lower bound on the complexity of implementing the driving Hamiltonian. To obtain overall accuracy $O(\varepsilon_{\text{tot}})$, the driving Hamiltonian needs accuracy of $O(\varepsilon_{\text{tot}}/\|H\|T)$ in the time. There are $\Theta(\|H\|T^2/\varepsilon_{\text{tot}})$ starting and finishing times, so by a counting argument the gate complexity is $\Omega(\log(\|H\|T/\varepsilon_{\text{tot}}))$. If the driving Hamiltonian is constant, then it is only the length of the time which is important, and that is limited to $O(1)$. The number of times is then $\Theta(\|H\|T/\varepsilon_{\text{tot}})$, but the lower bound on the complexity is still $\Omega(\log(\|H\|T/\varepsilon_{\text{tot}}))$. For constant error we therefore have $G = \Omega(\log(\|H\|T))$, as used in Section 3.

5 Measurement of the Control Qubits

What remains is to perform the final measurement. This should logically correspond to what happens if the state is decoded from its succinct representation to m qubits and then, for each qubit, an R gate is applied and it is measured in the computational basis. Of course, this cannot be literally implemented this way, because it would increase the gate and space usage to at least m ; our task is to *logically* perform this while remaining in the succinct representation.

Recall now that in Section 4.3, we constructed a procedure that approximately prepares $R^{\otimes n}|0^n\rangle$ (for any $n \leq m$) in succinct form [see $|\Xi_n^k\rangle$ in Eq. (8)]. We define U_n to be the ideal unitary that would exactly prepare the state $|\Xi_n^k\rangle$. The action of the ideal state preparation procedure is then $U_n C_n^k |0^n\rangle \approx C_n^k R^{\otimes n} |0^n\rangle$. The procedure we have described does not exactly perform this unitary, but it is within distance $O(\varepsilon)$. Also, we do not have an exact equality, because representations of terms with Hamming weight greater than k in $R^{\otimes n}|0^n\rangle$ are not obtained with the correct weights. More precisely, we have

$$U_n C_n^k |0^n\rangle = \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} C_n^k |x\rangle + \mu |\nu'\rangle. \quad (27)$$

This is to be compared with the uncompressed setting [Eq. (11)], in which we have

$$R^{\otimes n} |0^n\rangle = \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k}} \alpha^{n-|x|} \beta^{|x|} |x\rangle + \mu |\nu\rangle. \quad (28)$$

In terms of the *logical* data, U_n and $R^{\otimes n}$ produce almost the same state when applied to $|0^n\rangle$.

Returning to the issue of measurement, recall that initially we have prepared the state in a compressed form of $R^{\otimes m}|0^m\rangle$ (i.e., with $n = m$). In the *uncompressed* basis we would like to perform $R^{\otimes m}$, then perform a computational basis measurement. In the particular case that the computational basis measurement yielded all zeros, the measurement operator is $|0^m\rangle\langle 0^m| R^{\otimes m}$. Because we are performing all operations in the compressed basis, this measurement operator can be represented by $C_m^k |0^m\rangle\langle 0^m| R^{\otimes m} (C_m^k)^\dagger$. Because R is self-inverse, this is approximately the same as $C_m^k |0^m\rangle\langle 0^m| (C_m^k)^\dagger U_m^\dagger$. That is, to achieve this measurement result we first invert the preparation procedure described by U_m . Then, because $C_m^k |0^m\rangle = |m\rangle^{\otimes k+1}$ is a computational basis state, we can achieve the desired result by performing a computational basis measurement.

Ideally, this is what we want, but we also need to be able to find the positions of the ones in the case that the all-zero string is not obtained. At first glance, one might imagine that applying U_m^\dagger in place of $R^{\otimes m}$ would yield a succinct representation of the final outcome state, so measuring in the computational basis would provide the correct result. Unfortunately, this does not accurately simulate the final measurement except in the case where the all-zero string is obtained. The problem is that U_m and $R^{\otimes m}$ are *only* in close agreement when applied to the logical state $|0^m\rangle$. For any other logical state $|x\rangle$ (for non-zero $x \in \{0,1\}^m$), applying U_m and $R^{\otimes m}$ need not yield states in any close agreement.

Our first observation towards overcoming this problem is that we can at least perform an *incomplete* measurement that captures a seemingly small part of what we are seeking: we can cause the state to either collapse to logical $|0^m\rangle$ or to the subspace that is the orthogonal complement of this state—and with the correct probabilities. This is achieved by performing U_m^\dagger and then the 2-outcome incomplete projective measurement that distinguishes between the logical state $|0^m\rangle$ and its orthogonal complement $|0^m\rangle^\perp$, and then applying U_m to the resulting collapsed state. Our method to complete the measurement is to apply the above procedure recursively, on the two halves of the logical string. We now first motivate this procedure intuitively, followed by further technical details and a rigorous proof of correctness.

5.1 Measuring in Succinct Form: Intuition

The intuition behind our measurement strategy is given by the following simple thought experiment. Consider the problem of measuring an m -qubit state $|\psi\rangle$ in the computational basis. This can be accomplished by performing a sequence of two-outcome measurements in a variety of ways. One obvious approach is to measure the state of the first qubit, then the second qubit, and so on. Each final outcome $x \in \{0,1\}^m$ will occur with exactly the same probability as with the original complete measurement. We now describe an alternative—and unconventional—approach for simulating the same measurement.

First, perform the measurement distinguishing between $|0^m\rangle$ and $|0^m\rangle^\perp$, its orthogonal complement. If the state collapses to $|0^m\rangle$ we halt, outputting 0^m . Otherwise (when the state collapses to $|0^m\rangle^\perp$), apply the measurement $|0^{m/2}\rangle$ vs. $|0^{m/2}\rangle^\perp$ to the first $m/2$ qubits. If that part of the state collapses to $|0^{m/2}\rangle$ then output $0^{m/2}$ for the first $m/2$ bits; otherwise recurse further. Once this recursive measurement procedure for the first $m/2$ qubits has terminated, repeat it for the second $m/2$ qubits. Each final outcome $x \in \{0,1\}^m$ occurs with exactly the same probability as with the original complete measurement. Note that although this process may appear complicated, it terminates fast whenever the Hamming weight of the final outcome x is small: for Hamming weight up to k' , at most $k' \log m$ steps are performed.

Our actual scenario is different than the one described above in that the final measurement is in the basis $\{R^{\otimes m}|x\rangle : m \in \{0,1\}^m\}$ rather than the computational basis. However, our logical U_m and U_m^\dagger permit us to approximate the $R^{\otimes m}|0^m\rangle$ vs. $R^{\otimes m}|0^m\rangle^\perp$ measurement well. Also, making use of the fact that the underlying operation that we are simulating has a tensor product structure, $R^{\otimes m}|x_1x_2\rangle = R^{\otimes m/2}|x_1\rangle R^{\otimes m/2}|x_2\rangle$ for any $x_1, x_2 \in \{0,1\}^{m/2}$, we can emulate the recursive procedure in the above thought experiment. We now make this rigorous.

5.2 Measuring in Succinct Form: Details

We now introduce Alg. 4, which formalises the intuition behind the recursive measurement outlined above, and show that it simulates the desired measurement in succinct form. Recall that we assume without loss of generality that m is a power of 2.

Before stating Alg. 4, we require a lemma which allows us to efficiently “split” the encoded version of string $x = x_1x_2$ into the concatenation of the encoded versions of x_1 and x_2 .

Lemma 4 *Let $x = x_1x_2$ for $x \in \{0,1\}^n$ with $|x| \leq k$, $x_1, x_2 \in \{0,1\}^{n/2}$ and n a power of 2. Then there exists a quantum circuit with complexity $O(k \log n)$ for achieving the mapping*

$$C_n^k|x_1x_2\rangle \mapsto C_{n/2}^k|x_1\rangle \otimes C_{n/2}^k|x_2\rangle, \quad (29)$$

where $C_n^k|x_1x_2\rangle, C_{n/2}^k|x_1\rangle, C_{n/2}^k|x_2\rangle \in (\mathbb{C}^{n+1})^{\otimes k+1}$.

Proof. Because both $C_n^k|x_1x_2\rangle$ and $C_{n/2}^k|x_1\rangle \otimes C_{n/2}^k|x_2\rangle$ are computational basis states, the procedure that is performed is the same as would be performed classically, except that it must be performed coherently. That is, there is a reversible classical procedure to split the encoding in the computational basis, which immediately provides a coherent procedure for splitting the encoding. Because there are $O(k)$ registers of size $O(\log n)$, the complexity of this procedure is $O(k \log n)$ \square .

The formal statement of the recursive measurement algorithm is given in Alg. 4. To perform our recursive measurement, we simply call $\text{MEASURE}(A, 1, m)$, where A is the register

Algorithm 4 $S = \text{MEASURE}(A, m_1, m_2)$.

- Input: A – Registers corresponding to space $(\mathbb{C}^{m+1})^{\otimes k+1}$ containing the subset $\{m_1, \dots, m_2\}$ of the encoded control qubits.
 m_1 – The starting index $m_1 \in [m]$ of the encoded qubits in A .
 m_2 – The ending index $m_2 \in [m]$ of the encoded qubits in A .
- Precondition: $m_2 - m_1 + 1$ is a power of two.
- Output: A set of indices $S \subseteq [m]$ containing the positions where an uncompressed measurement would have found ones in the uncompressed setting.

Perform a measurement described by the measurement operators $M_{c,0}^{m_2-m_1+1}$ and $M_{c,1}^{m_2-m_1+1}$, where $M_{c,0}^n := U_n C_n^k |0^n\rangle\langle 0^n| (C_n^k)^\dagger U_n^\dagger$ and $M_{c,1}^n := I - M_{c,0}^n$. Label the measurement result d . Then

1. (Zero detected) If $d = 0$: Return $S = \emptyset$.
2. (Base case) If $d = 1$ and $m_1 = m_2$: Return $S = \{m_1\}$.
3. (Recurse) If $d = 1$ and $m_2 > m_1$: Split A to A_1 and A_2 , containing the encoded forms of the first and second halves, respectively, of the control qubits. Then return

$$S = \text{MEASURE}(A_1, m_1, (m_1+m_2-1)/2) \cup \text{MEASURE}(A_2, (m_1+m_2+1)/2, m_2). \quad (30)$$

containing our compressed control qubits. Once the procedure finishes running, it will return the locations of all the ones an uncompressed measurement would have obtained when measuring the uncompressed version of A . We truncate the recursive measurement procedure if k' ones have been located, to limit the complexity of the procedure.

We now introduce a notation that will be used throughout the remainder of the paper in order to simplify reference to quantities in the uncompressed protocol versus the compressed protocol. For quantities (states, operators or probabilities) in the compressed protocol, we will use a superscript or subscript “c”, whereas we will use “u” for the uncompressed protocol. To refer to quantities defined for both, we will use “ η ”. We also use n to refer to operations acting on a compressed sub-portion of the string of length n (instead of m for the full string).

To perform the measurement described by the measurement operators $M_{c,d}^n$ in Alg. 4, we apply U_n^\dagger , perform the measurement that distinguishes the encoded all-zero state from all other states, then apply U_n . In this form it is clear why we need to perform the operation U_n after the measurement: it means that all states orthogonal to that corresponding to measurement result 0 are unchanged, because they are just acted upon by the identity. The final U_n operation is also included for the 0 measurement result for simplicity, but it is not needed. As these measurement operators are projections, they are the same as the positive operator-valued measure elements.

For simplicity we have described the measurement in terms of the exact compressed measurement operators $M_{c,d}^n$ via the unitaries U_n and U_n^\dagger . We do not perform these measurements exactly, but the results are within $O(\varepsilon)$. More specifically, denoting the actual measurement operator that is performed by $\widetilde{M}_{c,d}^n$, we have the following result.

Lemma 5 *Our technique of approximately performing $M_{c,d}^n$ given above results in a trace*

distance error bounded as

$$\left\| \widetilde{M}_{c,d}^n \rho(\widetilde{M}_{c,d}^n)^\dagger - M_{c,d}^n \rho(M_{c,d}^n)^\dagger \right\|_{\text{tr}} \leq O(\varepsilon \text{Tr}(\rho)). \quad (31)$$

Note that this Lemma bounds the error in terms of the norm of the initial state, rather than the final state after the measurement. This means that there may be large error in the *normalised* state for a measurement result with low probability of occurring. This Lemma follows from linearity of the errors, but for completeness we provide a proof in Appendix A.

To show that the algorithm correctly simulates the desired uncompressed measurement, we consider a similar recursive measurement on the uncompressed state. We show that, except for the imprecision due to approximating U_n and omitting high Hamming weight components, the low Hamming weight portions of the states in Eqs. (27) and (28) evolve identically. Moreover, this holds even if the control qubits are entangled with a target register, as is generally the case here.

In the uncompressed setting, the state of the control and target registers before the final measurement can be described as approximately

$$|\widetilde{\psi}_u\rangle := \sum_{\substack{x \in \{0,1\}^m \\ |x| \leq k}} \gamma_{x,0}^{(0)} |x\rangle |w_x\rangle, \quad (32)$$

where $|w_x\rangle$ describes the state of the target register where the queries Q are applied. Note that $|\widetilde{\psi}_u\rangle$ is unnormalised, as we have omitted the high Hamming weight component. Similarly, in the compressed setting, before the final measurement we approximately have the (unnormalised) state

$$|\widetilde{\psi}_c\rangle := \sum_{\substack{x \in \{0,1\}^m \\ |x| \leq k}} \gamma_{x,0}^{(0)} C_m^k |x\rangle |w_x\rangle, \quad (33)$$

where the states $|w_x\rangle$ coincide with those in the uncompressed case. The coefficients $\gamma_{x,0}^{(0)}$ are the same in each case, and are equal to $i^{|x|} \xi_x^m$. We use this notation for consistency with the coefficients for the intermediate states in Eqs. (35) and (36) below.

We consider a measurement in the uncompressed case that is the same as in Alg. 4. We show that the results obtained in the two cases are close, but there are two sources of error: (1) the error incurred due to the high Hamming weight component of the state, and (2) the error due to not implementing U_n exactly. First we discuss the *error-free* case, i.e. where (1) we omit the high Hamming weight component, and where (2) U_n is implemented exactly. We subsequently reintroduce both sources of error and analyse their impacts. In the error-free analysis, we show the following.

Theorem 5 (Error-free simulation) *Assume we are in the error-free setting defined above. Then, suppose that before the final measurement, the states of the uncompressed and compressed control and target qubits are given by Eqs. (32) and (33), respectively. Then, Alg. 4 exactly simulates the uncompressed $R^{\otimes m}$ measurement in the following sense:*

1. *After running Alg. 4, the probability of obtaining a given measurement result is the same as for the uncompressed $R^{\otimes m}$ measurement, and*
2. *for a given measurement result the state of the target register in both uncompressed and compressed settings matches.*

Proof. Measuring $R^{\otimes m}|\tilde{\psi}\rangle$ in the computational basis can also be simulated using a recursive approach; namely, we apply $R^{\otimes m}$, followed by the incomplete measurement of $|0^m\rangle$ versus its orthogonal complement, then apply $R^{\otimes m}$. This can be represented by the measurement operators $M_{u,d}^m$, with

$$M_{u,0}^n := R^{\otimes n}|0^n\rangle\langle 0^n|R^{\otimes n}, \quad (34)$$

and $M_{u,1}^n := I - M_{u,0}^n$. Similar to Alg. 4, we are including the application of $R^{\otimes m}$ for both measurement results for simplicity, though it is not needed for result 0. If we obtain 1 as the outcome, we recurse on the two blocks of $m/2$ qubits by applying the measurement with operators $M_{u,d}^{m/2}$, and so forth.

To prove the result, we simply need to show that at each step in the recursion the states resulting from measurement operators $M_{c,d}^n$ and $M_{u,d}^n$ are equivalent. Let us denote the measurement result obtained at each step in the recursive measurement scheme by d_j . Then, at step ℓ , we have measurement results $d_1, \dots, d_{\ell-1}$, and will have a state that depends on those measurement results. Let us assume that at this step we have equivalent states for the compressed and uncompressed cases. The base case is that for $\ell = 1$, where the initial states (32) and (33) are equivalent. Then the states for the two cases can be expressed as

$$|\psi_{c,\ell-1}^{(d_1,\dots,d_{\ell-1})}\rangle = \sum_{|x|\leq k} \gamma_{x,\ell-1}^{(d_1,\dots,d_{\ell-1})} (C_n^k \otimes C_{\text{rest}}^k) |x\rangle |w_x\rangle, \quad (35)$$

$$|\psi_{u,\ell-1}^{(d_1,\dots,d_{\ell-1})}\rangle = \sum_{|x|\leq k} \gamma_{x,\ell-1}^{(d_1,\dots,d_{\ell-1})} |x\rangle |w_x\rangle. \quad (36)$$

At this stage the encoding will be a succinct encoding on a subset of n of the digits of x , and another encoding of the remaining digits (denoted C_{rest}^k), the exact form of which is unimportant for this analysis. The subset of n of the digits of x will depend on $d_1, \dots, d_{\ell-1}$. This dependence has not been indicated here for brevity. We also omit $x \in \{0,1\}^m$ from the sum for brevity.

In order for the results obtained for the compressed and uncompressed cases to be equivalent, all that is required is that the amplitude weightings $\gamma_{x,\ell-1}^{(d_1,\dots,d_{\ell-1})}$ in Eqs. (35) and (36) are the same. The results are equivalent in the sense that the probability of the measurement results, as well as the state of the target system for a given measurement result, are the same. The probability of the measurement results will be obtained from the normalisation of the state, which must be the same if the amplitudes are the same. Similarly the resulting state in the target system will be the same if the amplitudes are the same.

We will adopt the notation that I_{rest} indicates the identity on the remaining registers, so the overall measurement operator is $M_{c,d}^n \otimes I_{\text{rest}}$. We will also adopt the notation that x_n is the subset of n digits of the string x , and x_{rest} is the remaining digits. Then we have, using Lemma 1,

$$\langle 0^n | R^{\otimes n} | x_n \rangle = \langle x_n | R^{\otimes n} | 0^n \rangle = \langle x_n | (C_n^k)^\dagger U_n C_n^k | 0^n \rangle = \langle 0^n | (C_n^k)^\dagger U_n^\dagger C_n^k | x_n \rangle. \quad (37)$$

For the compressed case, consider performing the measurement with operators $M_{c,d}^n$. In the

case that the measurement result is $d = 0$, our compressed state becomes

$$\begin{aligned}
& (M_{c,0}^n \otimes I_{\text{rest}}) |\psi_{c,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle \\
& \approx (U_n C_n^k |0^n\rangle \langle 0^n| (C_n^k)^\dagger U_n^\dagger \otimes I_{\text{rest}}) \left[\sum_{|x| \leq k} \gamma_{x,\ell-1}^{(d_1, \dots, d_{\ell-1})} (C_n^k \otimes C_{\text{rest}}^k) |x\rangle |w_x\rangle \right] \\
& = \sum_{|x| \leq k} \gamma_{x,\ell-1}^{(d_1, \dots, d_{\ell-1})} (\langle 0^n | (C_n^k)^\dagger U_n^\dagger C_n^k |x_n\rangle) U_n C_n^k |0^n\rangle C_{\text{rest}}^k |x_{\text{rest}}\rangle |w_x\rangle \\
& = \sum_{|x| \leq k} \gamma_{x,\ell-1}^{(d_1, \dots, d_{\ell-1})} \langle 0^n | R^{\otimes n} |x_n\rangle U_n C_n^k |0^n\rangle C_{\text{rest}}^k |x_{\text{rest}}\rangle |w_x\rangle \\
& \approx \sum_{|x| \leq k} \gamma_{x,\ell-1}^{(d_1, \dots, d_{\ell-1})} \langle 0^n | R^{\otimes n} |x_n\rangle \left(\sum_{|y| \leq k} \xi_y^n C_n^k |y\rangle \right) C_{\text{rest}}^k |x_{\text{rest}}\rangle |w_x\rangle =: |\tilde{\psi}_{c,\ell}^{(d_1, \dots, d_{\ell-1}, 0)}\rangle, \quad (38)
\end{aligned}$$

where $\xi_y^n = \alpha^{n-|y|} \beta^{|y|}$, and y is an n -digit string. The approximate equality in the first line of Eq. (38) is because the measurement operator $M_{c,0}^n$ cannot be obtained exactly, because the unitary U_n is not performed exactly. The approximate equality in the last line is because the high Hamming weight components have been omitted. In the error-free setting the error in these approximations is ignored.

In comparison, in the uncompressed setting, a similar calculation yields, for $d = 0$,

$$\begin{aligned}
& (M_{u,0}^n \otimes I_{\text{rest}}) |\psi_{u,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle = (R^{\otimes n} |0^n\rangle \langle 0^n| R^{\otimes n} \otimes I_{\text{rest}}) |\psi_{u,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle \\
& \approx \sum_{|x| \leq k} \gamma_{x,\ell-1}^{(d_1, \dots, d_{\ell-1})} \langle 0^n | R^{\otimes n} |x_n\rangle \left(\sum_{|y| \leq k} \xi_y^n |y\rangle \right) |x_{\text{rest}}\rangle |w_x\rangle =: |\psi_{u,\ell}^{(d_1, \dots, d_{\ell-1}, 0)}\rangle. \quad (39)
\end{aligned}$$

The approximate equality in the last line is again due to omitting high Hamming weight components. In the error-free setting the error in this approximation is ignored. In the case that the measurement result is $d = 1$, then the states obtained are

$$\begin{aligned}
& (I - M_{c,0}^n \otimes I_{\text{rest}}) |\psi_{c,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle = |\psi_{c,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle - |\tilde{\psi}_{c,\ell}^{(d_1, \dots, d_{\ell-1}, 0)}\rangle =: |\tilde{\psi}_{c,\ell}^{(d_1, \dots, d_{\ell-1}, 1)}\rangle, \\
& (I - M_{u,0}^n \otimes I_{\text{rest}}) |\psi_{u,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle = |\psi_{u,\ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle - |\psi_{u,\ell}^{(d_1, \dots, d_{\ell-1}, 0)}\rangle =: |\psi_{u,\ell}^{(d_1, \dots, d_{\ell-1}, 1)}\rangle. \quad (40)
\end{aligned}$$

Above we have defined resulting states after the measurements in the uncompressed and compressed setting of $|\psi_{u,\ell}^{(d_1, \dots, d_\ell)}\rangle$ and $|\tilde{\psi}_{c,\ell}^{(d_1, \dots, d_\ell)}\rangle$, respectively. The quantity $|\tilde{\psi}_{c,\ell}^{(d_1, \dots, d_\ell)}\rangle$ is the state in the compressed case before the change in the compression. To obtain the state $|\psi_{c,\ell}^{(d_1, \dots, d_\ell)}\rangle$, the compression of the string must be changed as per Lemma 4. This can be done without error, and does not change the amplitudes.

Omitting the high Hamming weight states, we start with states $|\tilde{\psi}_\eta\rangle$, which have the same amplitudes in the compressed and uncompressed cases. Then, by the above reasoning, if the amplitudes are the same at step $\ell - 1$, they are the same at step ℓ . Therefore, by induction, the amplitudes must be the same after the full recursive measurement. Therefore the same amplitudes are obtained for the compressed and uncompressed cases, so the results obtained in the compressed and uncompressed cases are equivalent. That is, the probabilities of the measurement results and the state of the target register for a given measurement result match \square .

Theorem 5 shows that if we focus solely on the low Hamming weight subspace, and if we assume we can prepare the state $C_n^k|0^n\rangle$ exactly, then our succinct recursive measurement Alg. 4 *perfectly* simulates the uncompressed measurement. We now analyse the error incurred when these two assumptions are dropped. First we need to identify the appropriate measure of the error in the measurement. We would like to bound the average trace distance; i.e.

$$\overline{D} := \sum_{\mathbf{b}} p_{\mathbf{b}}^c \|\rho_{\mathbf{b}}^u - \rho_{\mathbf{b}}^c\|_{\text{tr}}, \quad (41)$$

where $p_{\mathbf{b}}^u$ is the probability of obtaining the measurement result $\mathbf{b} = (b_1, \dots, b_m)$, and $\rho_{\mathbf{b}}^u$ is the state for the target system. We would also like to bound the error in the probabilities obtained. This is because measurement results with many ones will be difficult to correct, so we need to ensure that the probabilities for those measurement results remain small. The error in the probability distribution can be quantified by

$$\Delta p := \sum_{\mathbf{b}} |p_{\mathbf{b}}^u - p_{\mathbf{b}}^c|. \quad (42)$$

We can bound both those errors using the quantity

$$D_{\text{av}} := \sum_{\mathbf{b}} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}}. \quad (43)$$

Because the trace distance is non-increasing under channels, and we obtain Δp by applying the completely depolarising channel to both $\rho_{\mathbf{b}}^u$ and $\rho_{\mathbf{b}}^c$ in Eq. (43), we have $\Delta p \leq D_{\text{av}}$. Then we have

$$p_{\mathbf{b}}^c \|\rho_{\mathbf{b}}^u - \rho_{\mathbf{b}}^c\|_{\text{tr}} \leq \|p_{\mathbf{b}}^c \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u\|_{\text{tr}} + \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}} \leq 2 \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}}. \quad (44)$$

Summing over \mathbf{b} then gives $\overline{D} \leq 2D_{\text{av}}$.

Theorem 6 (Error bounds) *The error between compressed and uncompressed schemes can be bounded as*

$$D_{\text{av}} = O(\varepsilon' + \varepsilon k' \log m). \quad (45)$$

Proof. In order to bound the value of D_{av} , we have four main sources of error:

1. Omitting the high Hamming weight components of the initial states.
2. Omitting measurement results with Hamming weight greater than k' .
3. Omitting high Hamming weight components in each step of the recursive measurement.
4. Inaccuracy in performing the U_n operations in each step of the recursive measurement.

Error sources 1 and 2 introduce error $O(\varepsilon)$ and $O(\varepsilon')$, respectively. The contribution of the error from sources 3 and 4 may be bounded as follows. In locating the position of a single one in the measurement result, there is a contribution of $O(\varepsilon)$ to the error from each of the steps as described in Eq. (38). These need to be performed $\log m$ times, and as a result the contribution to the error is $O(\varepsilon \log m)$. If h ones need to be located, the worst case is where the sequence of measurements to locate these ones is independent, so the total contribution

to the error from sources 3 and 4 is $O(h\varepsilon \log m)$. Since the error due to locating no more than k' ones will be $O(\varepsilon')$, we can take $h \leq k'$, and bound the overall error by $O(\varepsilon k' \log m + \varepsilon')$. A more rigorous form of this proof is given in Appendix B \square .

Error sources 3 and 4 give a contribution to the error of $O(\varepsilon)$ times the norm of the state for each step of the recursive measurement. However, for many initial sequences of measurement results, at step ℓ all ones have already been located, so there are no further measurements needed. This means that the measurements at this point are just the identity, and no further error is introduced for that sequence of initial measurement results. This means that bounding the additional error by $O(\varepsilon)$ for each ℓ overestimates the error. We will show that the error can be bound by using the mean number of ones that are measured. In the case of the uncompressed measurements, the probability of each one is $\leq 2\alpha^2\beta^2$. Because $\beta^2 \approx 1/8m$, the expected number of ones is $\leq 2\beta^2m = O(1)$.

Theorem 7 (Improved error bounds) *Provided $\varepsilon = O(1/(k' \log m))$, the error between the compressed and uncompressed schemes can be bounded as*

$$D_{\text{av}} = O(\varepsilon' + \varepsilon \log m). \quad (46)$$

Proof. More specifically, $\rho_{\eta, \ell-1}$ (defined in Eq. (B.15)) will have a component where the recursive measurement scheme has not terminated yet, and another measurement needs to be performed. This component will be that where the ancillas contain $d_1, \dots, d_{\ell-1}$ corresponding to sequences of measurement results such that further measurements need to be performed. There will also be a component corresponding to sequences of measurement results where the recursive measurement scheme has finished. We will denote the components corresponding to that where the recursive measurement has not terminated or has terminated by $\rho_{\eta, \ell-1}^{\text{con}}$ and $\rho_{\eta, \ell-1}^{\text{fin}}$, respectively. More explicitly, if we denote by S_{con} and S_{fin} the sets of measurement results $(d_1, \dots, d_{\ell-1})$ that correspond to a recursive measurement that has not terminated or has terminated, respectively, then we have

$$\rho_{\eta, \ell-1}^{\text{con/fin}} := \sum_{(d_1, \dots, d_{\ell-1}) \in S_{\text{con/fin}}} |d_{\ell-1}\rangle \langle d_{\ell-1}| \otimes \dots \otimes |d_1\rangle \langle d_1| \otimes |\psi_{\eta, \ell-1}^{(d_1, \dots, d_{\ell-1})}\rangle \langle \psi_{\eta, \ell-1}^{(d_1, \dots, d_{\ell-1})}|. \quad (47)$$

Because the measurement acts only on $\rho_{\eta, \ell-1}^{\text{con}}$, and the error in the measurement is bound by $O(\varepsilon)$ times the trace of the state the measurement acts upon, the error in approximating $\mathcal{E}_{\eta, \ell}(\rho_{\eta, \ell-1})$ by $\rho_{\eta, \ell}$ ($\mathcal{E}_{\eta, \ell}$ is defined in Eq. (B.12)) will be bounded by $O(\varepsilon \text{Tr}(\rho_{\eta, \ell-1}^{\text{con}}))$. Therefore the total error from sources 3 and 4 is bounded by

$$O\left(\varepsilon \sum_{\ell=1}^K \text{Tr}(\rho_{\eta, \ell-1}^{\text{con}})\right). \quad (48)$$

But, because the number of measurement steps need be no larger than $1 + h \log m$, where h is the number of ones found by the measurement, the probability that the number of ones is $\geq h$ is no less than $\text{Tr}(\rho_{\eta, h \log m}^{\text{con}})$. Denoting the probability that the number of ones is $\geq h$ by

$p(|\mathbf{b}| \geq h)$, we can bound the sum of the traces by

$$\begin{aligned} \sum_{\ell=1}^K \text{Tr}(\rho_{\eta, \ell-1}^{\text{con}}) &\leq \sum_{\ell=1}^K \text{Tr}(\rho_{\eta, \lfloor (\ell-1)/\log m \rfloor \log m}^{\text{con}}) \\ &\leq \log m \sum_{h=0}^m p(|\mathbf{b}| \geq h) = (\langle |\mathbf{b}| \rangle + 1) \log m. \end{aligned} \quad (49)$$

Next we need to take into account the fact that here the expectation value of the number of ones is for the approximate states $\rho_{u, \ell-1}^{\text{con}}$, not for the exact uncompressed measurement scheme. To take account of this difference, we can use the cumulative error to bound the error in the norm of the state at each step. Note that the norm of $\rho_{\eta, \ell-1}^{\text{con}}$ is the same for $\eta = u$ and $\eta = c$, so we only need perform the analysis for $\eta = u$. Let $A_{\ell-1}$ denote the norm, for the exact uncompressed measurement, of the component where the recursive measurement scheme has not stopped before step ℓ . In addition, let $E_{\ell-1}$ denote the cumulative error before step ℓ . Then the increment in the error is bound by ε times the norm of the non-terminated component, which is bound by $A_{\ell-1}$ plus the cumulative error.

$$\begin{aligned} E_{\ell} &\leq E_{\ell-1} + O(\varepsilon A_{\ell-1} + \varepsilon E_{\ell-1}) \\ &= E_{\ell-1}[1 + O(\varepsilon)] + O(\varepsilon A_{\ell-1}). \end{aligned} \quad (50)$$

Multiplying both sides by $[1 + O(\varepsilon)]^{K-\ell}$, we obtain

$$E_{\ell}[1 + O(\varepsilon)]^{K-\ell} \leq E_{\ell-1}[1 + O(\varepsilon)]^{K-(\ell-1)} + O(\varepsilon A_{\ell-1}[1 + O(\varepsilon)]^{K-\ell}). \quad (51)$$

As a result, the final error is bound by

$$\begin{aligned} E_K &\leq \sum_{\ell=1}^K O(\varepsilon A_{\ell-1}[1 + O(\varepsilon)]^{K-\ell}) \\ &\leq [1 + O(\varepsilon)]^K \sum_{\ell=1}^K O(\varepsilon A_{\ell-1}) \\ &\leq O(\exp(\varepsilon K) \varepsilon (\langle |\mathbf{b}| \rangle + 1) \log m). \end{aligned} \quad (52)$$

Here the expectation value of the number of ones is for the exact scheme, which is $O(1)$, and we therefore find that the error is bound by $O(\exp(\varepsilon K) \varepsilon \log m)$. Recall that we take $K = O(k' \log m)$. This means that, provided $\varepsilon = O(1/(k' \log m))$, $\exp(\varepsilon K)$ is $O(1)$, and we obtain scaling of the error of $O(\varepsilon \log m)$. Adding $O(\varepsilon' + \varepsilon)$ to take account of error sources 1 and 2 yields the result given in the Theorem \square .

Given the conditions of this Theorem, the overall error for each time step is $O(\varepsilon' + \varepsilon \log m)$. This includes error in simulating the driving Hamiltonian. The driving Hamiltonian may be applied up to k' times, though the expected number of times is $O(1)$. As the allowable error in the driving Hamiltonian is $O(\varepsilon')$, that gives a contribution of $O(\varepsilon')$ to the error in each time step. As there are $O(T)$ time steps, the total error is $O(\varepsilon' T + \varepsilon T \log m)$. To limit the error of the overall scheme to ε_{tot} , we take $\varepsilon' = O(\varepsilon_{\text{tot}}/T)$ and $\varepsilon = O(\varepsilon_{\text{tot}}/(T \log m))$. Then $k' = O(\log(1/\varepsilon')) = O(\log(T/\varepsilon_{\text{tot}}))$. As we consider large T and small ε_{tot} , we therefore have $\varepsilon = O(1/[\log(T/\varepsilon_{\text{tot}}) \log m]) = O(1/(k' \log m))$. This means that the condition of the Theorem is satisfied with this choice of parameters, and the total error will be bounded by ε_{tot} .

6 Proof of Main Theorem

Finally we are in a position to prove Theorem 1.

Proof. (of Theorem 1) First, the number of oracle queries is $O(k'T)$, because we have divided the simulation into $O(T)$ time intervals, and limit the number of queries required within each time interval to $O(k')$. The value of k' is chosen to ensure that the error due to omitting high Hamming-weight states $O(1)$ times within each time interval is no more than ε' . We can bound the total error by ε_{tot} if we take $\varepsilon' = O(\varepsilon_{\text{tot}}/T)$, which means that k' scales as

$$k' = O\left(\frac{\log(T/\varepsilon_{\text{tot}})}{\log \log(T/\varepsilon_{\text{tot}})}\right). \quad (53)$$

Then the overall number of oracle calls scales as

$$O\left(\frac{T \log(T/\varepsilon_{\text{tot}})}{\log \log(T/\varepsilon_{\text{tot}})}\right). \quad (54)$$

Omitting the dependence on ε_{tot} gives the result given in the statement of the Theorem.

Next we discuss the number of gates required for Alg. 4. The maximum number of steps in the recursive procedure is $1 + 2k' \log m$, but the expected number of steps is $O(\log m)$. For the full algorithm for the evolution over time T , there are many of these recursive measurements, and the probability of the average number of steps differing significantly from its expected value is small. Similarly to the analysis in Section 4.1, an upper bound of $O(1/\varepsilon_{\text{tot}})$ times the average value will not be exceeded with probability $1 - O(\varepsilon_{\text{tot}})$. As ε_{tot} is taken to be constant, this does not affect the final result. Because U_n and U_n^\dagger are performed at each step, these operations are performed $O(\log m)$ times. As was found above, the complexity of the operation U_n is $O(k[\log m + \log \log(1/\varepsilon)])$. Therefore the overall complexity for this time step is $O(k[(\log m)^2 + \log m \log \log(1/\varepsilon)])$.

It is also necessary to perform $O(k')$ time evolutions under the driving Hamiltonian. In the definition of the problem we let G be the number of gates required for the simulation of the driving Hamiltonian, so that the number of gates to simulate the driving Hamiltonian in this time step is $O(k'G)$. Therefore, the scaling for the total number of gates is

$$O(TGk' + Tk[(\log m)^2 + \log m \log \log(1/\varepsilon)]). \quad (55)$$

Next we use $\varepsilon' = O(\varepsilon_{\text{tot}}/T)$ and $\varepsilon = O(\varepsilon_{\text{tot}}/(T \log m))$. As discussed in Section 4.6, we can take $\log m = O(\log(\|H\|T/\varepsilon_{\text{tot}}))$. Considering the scaling with large $\|H\|$, the total number of gates simplifies to

$$O\left(\frac{TG \log(T/\varepsilon_{\text{tot}})}{\log \log(T/\varepsilon_{\text{tot}})} + \frac{T \log[(T \log m)/\varepsilon_{\text{tot}}]}{\log \log[(T \log m)/\varepsilon_{\text{tot}}]} (\log m)^2\right). \quad (56)$$

A further simplification may be obtained by ignoring the double-log factors in the denominators, and then using the scaling of $\log m$ to give

$$O(TG \log(T/\varepsilon_{\text{tot}}) + T[\log(T/\varepsilon_{\text{tot}}) + \log \log \|H\|][\log(T/\varepsilon_{\text{tot}}) + \log \|H\|]^2). \quad (57)$$

The number of gates can then be bounded in a simpler but looser form as

$$O(TG \log(T/\varepsilon_{\text{tot}}) + T[\log(\|H\|T/\varepsilon_{\text{tot}})]^3). \quad (58)$$

Omitting ε_{tot} , because we take this quantity to be constant, gives the scaling in the Theorem.

The number of qubits required for the algorithm is dominated by the number of qubits required for the recursive measurement scheme. The number of qubits used for the ancilla space is $O(k[\log m + \log \log(1/\varepsilon)])$. In the recursive measurement scheme it may be necessary to duplicate the ancilla space k' times to ensure that a maximum of k' ones are detected. The overall space used is therefore

$$O\left(\frac{\log[(T \log m)/\varepsilon_{\text{tot}}]}{\log \log[(T \log m)/\varepsilon_{\text{tot}}]} \frac{\log(T/\varepsilon_{\text{tot}})}{\log \log(T/\varepsilon_{\text{tot}})} [\log m + \log \log(T \log m/\varepsilon_{\text{tot}})]\right). \quad (59)$$

Cancelling the double-log, then omitting double-log factors in the denominator gives

$$O(\log[(T \log m)/\varepsilon_{\text{tot}}] \log(T/\varepsilon_{\text{tot}}) \log m). \quad (60)$$

Using the scaling of m then gives

$$O(\log(T/\varepsilon_{\text{tot}})[\log(T/\varepsilon_{\text{tot}}) + \log \log \|H\|] \log(\|H\|T/\varepsilon_{\text{tot}})). \quad (61)$$

A simpler bound can be given as

$$O([\log(\|H\|T/\varepsilon_{\text{tot}})]^3). \quad (62)$$

Again omitting ε_{tot} gives the scaling in the statement of the Theorem.

Note also that the allowable error in the driving Hamiltonian is $O(\varepsilon')$, which is $O(\varepsilon_{\text{tot}}/T)$. For constant ε_{tot} , the allowable error in the implementation of the driving Hamiltonian is $O(1/T)$, as given in the statement of the Theorem \square .

7 Conclusions

We have shown that any continuous-time query algorithm of cost T can be implemented with a number of discrete queries close to linear in T , and with a number of gates that is also close to linear in T . This means that any continuous-time quantum algorithm can be converted into an efficient discrete-query algorithm. In contrast, using the algorithm of Ref. [2] directly would result in a number of gates that is linear in mT . That is, the gate complexity would be superlinear in $\|H\|T$, and similar to what would be obtained just using product formulae.

Our results provide an even better improvement in the scaling with $\|H\|$; the number of gates is polylogarithmic in this quantity, rather than superlinear. As the norm of the driving Hamiltonian can potentially be very large, this can potentially provide a very large improvement in efficiency. In both cases, the query complexity is independent of $\|H\|$, but it does not appear to be possible to completely remove the dependence of the number of gates on $\|H\|$ via this approach.

The methods we have presented may also be used as an alternative to product formulae when simulating state evolution for a sum of Hamiltonians, where one Hamiltonian is self-inverse, and the other has large norm, $\|H\|$. Previous work has considered the complexity of Hamiltonian simulation via product formulae where one Hamiltonian has much larger norm [11]. Even using that approach, the complexity is only reduced from $O(\|H\|T(\|H\|T/\varepsilon_{\text{tot}})^\delta)$ to $O(\|H\|T(T/\varepsilon_{\text{tot}})^\delta)$. In comparison, here we have obtained complexity that is polylogarithmic in $\|H\|$.

Acknowledgements

Research supported by Canada's NSERC, CIFAR, MITACS, the U.S. ARO, and ARC grant FT100100761.

References

1. E. Farhi and S. Gutmann (1998), *Analog analogue of a digital quantum computation*, Phys. Rev. A, 57, pp. 2403–2406.
2. R. Cleve, D. Gottesman, M. Mosca, R. Somma, and D. Yonge-Mallo (2009), *Efficient discrete-time simulations of continuous-time quantum query algorithms*, In *Proc. 41st ACM Symposium on Theory of Computing*, pp. 409–416.
3. T. Lee, R. Mittal, B. W. Reichardt, R. Špalek, and M. Szegedy (2011), *Quantum query complexity of state conversion*, In *Proc. 52nd IEEE Symposium on Foundations of Computer Science*, pp. 344–353; arXiv:1011.3020.
4. E. Farhi, J. Goldstone, and S. Gutmann (2008), *A quantum algorithm for the Hamiltonian NAND tree*, Theory of Computing, 4, pp. 169–190.
5. A. M. Childs, R. Cleve, S. P. Jordan, and D. Yonge-Mallo (2009), *Discrete-query quantum algorithm for NAND trees*, Theory of Computing, 5, pp. 119–123.
6. A. Ambainis, A. M. Childs, B. W. Reichardt, R. Špalek, and S. Zhang (2007), *Any AND-OR formula of size N can be evaluated in time $N^{1/2+o(1)}$ on a quantum computer*, In *Proc. 48th IEEE Symposium on Foundations of Computer Science*, pp. 363–372.
7. D. W. Berry, G. Ahokas, R. Cleve, and B. C. Sanders (2007), *Efficient quantum algorithms for simulating sparse Hamiltonians*, Commun. Math. Phys., 270, pp. 359–371.
8. L. Grover and T. Rudolph (2002), *Creating superpositions that correspond to efficiently integrable probability distributions*, arXiv:quant-ph/0208112.
9. N. Wiebe, D. W. Berry, P. Høyer, and B. C. Sanders (2010), *Higher order decompositions of ordered operator exponentials*, J. Phys. A: Math. Theor., 43, 065203.
10. J. Huyghebaert and H. De Raedt (1990), *Product formula methods for time-dependent Schrodinger problems*, J. Phys. A: Math. Gen., 23, pp. 5777–5793.
11. A. Papageorgiou and C. Zhang (2012), *On the efficiency of quantum algorithms for Hamiltonian simulation*, Quantum Information Processing, 11, pp. 541–561; arXiv:1005.1318.

Appendix A

The difference between the desired measurement operator $M_{c,d}^n$ and the actual measurement operator $\widetilde{M}_{c,d}^n$ is because we will use operations on an expanded space that includes an error-flag ancilla. Recall that, because $|\phi_{q-t}\rangle$ is not exactly equal to $|\phi_q\rangle$, we have a register that is not exactly reset to zero, and this is swapped into an ancilla register. The unitary operations in this expanded space will be denoted \widetilde{U}_n and \widetilde{U}_n^\dagger . Then the action of \widetilde{U}_n is

$$\widetilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle = \sum_{\substack{x \in \{0,1\}^n \\ |x| \leq k+1}} \xi_x^n [\sqrt{1-\varepsilon_x} C_n^k |x\rangle \otimes |0\rangle + \sqrt{\varepsilon_x} |\text{err}_x\rangle \otimes |1\rangle]. \quad (\text{A.1})$$

Here the tensor product with $|0\rangle$ on the left-hand side indicates the use of ancillas that are initially in the state zero. The amplitudes ξ_x^n are the amplitudes for each $C_n^k |x\rangle$ in the state (8). These amplitudes include those for $|x| = k+1$ for the state $|\nu'\rangle$, which corresponds to encoded Hamming-weight $k+1$ states. For $|x| \leq k$, we have $\xi_x^n = \alpha^{n-|x|} \beta^{|x|}$. The tensor product with $|0\rangle$ on the right-hand side indicates ancillas that will be set to zero in the case

of success. The parameter ε_x is $\leq \varepsilon$, and can in general depend on x . The state $|\text{err}_x\rangle$ is an error state.

For the ideal state preparation, we have

$$\langle x|(C_n^k)^\dagger U_n C_n^k |0^n\rangle = \xi_x^n. \quad (\text{A.2})$$

Using the expression for the action of \tilde{U}_n , we find

$$[(\langle x|(C_n^k)^\dagger) \otimes \langle 0|][\tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle] = \xi_x^n \sqrt{1 - \varepsilon_x} = \langle x|(C_n^k)^\dagger U_n C_n^k |0^n\rangle [1 - O(\varepsilon)]. \quad (\text{A.3})$$

There is no contribution from the error register, because the error flag is orthogonal to zero for that register.

To perform the measurement, we append ancillas in the zero state, and perform \tilde{U}_n^\dagger . Then we perform the measurement that projects onto $C_n^k |0^n\rangle \otimes |0\rangle$ and its orthogonal complement. Here the tensor product with $|0\rangle$ indicates the extra ancillas used by the full preparation procedure \tilde{U}_n . Then we perform \tilde{U}_n .

The action of this measurement will have error $O(\varepsilon)$ from that used in the algorithm. First, consider the resulting state for 0 measurement result and initial state $C_n^k |x\rangle$.

$$\begin{aligned} \tilde{M}_{c,0}^n C_n^k |x\rangle \otimes |0\rangle &= \tilde{U}_n [C_n^k |0^n\rangle \langle 0^n| (C_n^k)^\dagger \otimes |0\rangle \langle 0|] \tilde{U}_n^\dagger C_n^k |x\rangle \otimes |0\rangle \\ &= \tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle [\langle x| \langle 0| (C_n^k)^\dagger \tilde{U}_n C_n^k |0^n\rangle |0\rangle]^* \\ &= \tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle [\langle x| (C_n^k)^\dagger U_n C_n^k |0^n\rangle]^* [1 - O(\varepsilon)] \\ &= \tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle [\langle 0^n| (C_n^k)^\dagger U_n^\dagger C_n^k |x\rangle] [1 - O(\varepsilon)]. \end{aligned} \quad (\text{A.4})$$

The action of the exact measurement operator is

$$M_{c,0}^n C_n^k |x\rangle |0\rangle = U_n C_n^k |0^n\rangle \otimes |0\rangle [\langle 0^n| (C_n^k)^\dagger U_n^\dagger C_n^k |x\rangle]. \quad (\text{A.5})$$

In addition, $\tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle$ is an approximation of $U_n C_n^k |0^n\rangle \otimes |0\rangle$ with trace distance $O(\varepsilon)$. It is therefore clear that the trace distance between $\tilde{M}_{c,0}^n C_n^k |x\rangle |0\rangle$ and $M_{c,0}^n C_n^k |x\rangle |0\rangle$ is $O(\varepsilon)$. Similarly, the trace distance will be $O(\varepsilon)$ for any normalised pure state superposition of $|x\rangle$. By convexity of trace distance, for state ρ , the trace distance will be $O(\varepsilon \text{Tr}(\rho))$.

The resulting state for measurement result 1 is then

$$\tilde{M}_{c,1}^n C_n^k |x\rangle \otimes |0\rangle = C_n^k |x\rangle |0\rangle - \tilde{U}_n C_n^k |0^n\rangle \otimes |0\rangle [\langle 0| (C_n^k)^\dagger U_n^\dagger C_n^k |x\rangle] [1 - O(\varepsilon)]. \quad (\text{A.6})$$

This is because \tilde{U}_n is exactly the inverse of \tilde{U}_n^\dagger in the expanded space. The error for measurement result 1 is therefore the same as for measurement result 0. Therefore, for both measurement results the trace distance for initial state ρ will be $O(\varepsilon \text{Tr}(\rho))$.

Appendix B

To make the analysis of Theorem 6 rigorous, we first want to omit the high Hamming weight measurement results. For the measurements in the uncompressed case, the probability of measurement results with Hamming weight over k' is $O(\varepsilon')$. This is because the probability of obtaining each one is no more than $2\alpha^2\beta^2$. Because we take $\beta^2 \approx 1/8m$, the probability of obtaining more than k' ones with $k' = \Theta(\log(1/\varepsilon')/\log \log(1/\varepsilon'))$ is $O(\varepsilon')$. Recall that we

place a bound ε' on errors that only occur once in each time step, and use a corresponding Hamming weight cutoff k' , whereas we use k for limiting errors that occur multiple times in the measurement process.

To bound D_{av} , we also need to take account of the probability of high Hamming weight measurement results for the uncompressed measurement. We can do this in the following way. First use

$$\sum_{|\mathbf{b}| > k'} (p_{\mathbf{b}}^c - p_{\mathbf{b}}^u) = \sum_{|\mathbf{b}| \leq k'} (p_{\mathbf{b}}^u - p_{\mathbf{b}}^c) \leq \sum_{|\mathbf{b}| \leq k'} |p_{\mathbf{b}}^u - p_{\mathbf{b}}^c| \leq \sum_{|\mathbf{b}| \leq k'} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}}. \quad (\text{B.1})$$

Therefore we can bound D_{av} by

$$\begin{aligned} D_{\text{av}} &\leq \sum_{|\mathbf{b}| > k'} (p_{\mathbf{b}}^c + p_{\mathbf{b}}^u) + \sum_{|\mathbf{b}| \leq k'} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}} \\ &= \sum_{|\mathbf{b}| > k'} (p_{\mathbf{b}}^c - p_{\mathbf{b}}^u) + 2 \sum_{|\mathbf{b}| > k} p_{\mathbf{b}}^u + \sum_{|\mathbf{b}| \leq k'} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}} \\ &\leq O(\varepsilon') + 2 \sum_{|\mathbf{b}| \leq k'} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}}. \end{aligned} \quad (\text{B.2})$$

This means that omitting the high Hamming weight measurement results can only change the results by a multiplying factor and an $O(\varepsilon')$ term. For convenience we define

$$D'_{\text{av}} := \sum_{|\mathbf{b}| \leq k'} \|p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c\|_{\text{tr}}. \quad (\text{B.3})$$

Next we note that the distance measure can be written as a trace distance between two states, rather than the average of trace distances. That is,

$$D'_{\text{av}} = \left\| \sum_{|\mathbf{b}| \leq k'} (p_{\mathbf{b}}^u |\mathbf{b}\rangle \langle \mathbf{b}| \otimes \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c |\mathbf{b}\rangle \langle \mathbf{b}| \otimes \rho_{\mathbf{b}}^c) \right\|_{\text{tr}}. \quad (\text{B.4})$$

The reason for this is that the complete matrix is block-diagonal, with $p_{\mathbf{b}}^u \rho_{\mathbf{b}}^u - p_{\mathbf{b}}^c \rho_{\mathbf{b}}^c$ as the blocks on the diagonal. The trace distance for the entire density matrix is just the sum of the trace distances for the blocks on the diagonal, which is the definition of D'_{av} .

Let us denote by $|\psi_{\eta}\rangle$ the states obtained after preparation and controlled operations. Then we have

$$p_{\mathbf{b}}^{\eta} \rho_{\mathbf{b}}^{\eta} = \text{Tr}_{\text{ctrl}}(M_{\eta, \mathbf{b}} |\psi_{\eta}\rangle \langle \psi_{\eta}| M_{\eta, \mathbf{b}}^{\dagger}). \quad (\text{B.5})$$

Here Tr_{ctrl} indicates a trace over the control registers. Then we have

$$D'_{\text{av}} = \left\| \sum_{|\mathbf{b}| \leq k'} \left[|\mathbf{b}\rangle \langle \mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\mathbf{u}, \mathbf{b}} |\psi_{\mathbf{u}}\rangle \langle \psi_{\mathbf{u}}| M_{\mathbf{u}, \mathbf{b}}^{\dagger} \right) - |\mathbf{b}\rangle \langle \mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\mathbf{c}, \mathbf{b}} |\psi_{\mathbf{c}}\rangle \langle \psi_{\mathbf{c}}| M_{\mathbf{c}, \mathbf{b}}^{\dagger} \right) \right] \right\|_{\text{tr}}. \quad (\text{B.6})$$

Now note that the maps defined by

$$\mathcal{E}_{\eta}(\rho) := \sum_{\mathbf{b}} |\mathbf{b}\rangle \langle \mathbf{b}| \otimes \text{Tr}_{\text{ctrl}}(M_{\eta, \mathbf{b}} \rho M_{\eta, \mathbf{b}}^{\dagger}), \quad (\text{B.7})$$

are completely-positive trace-preserving (CPTP). This means that trace distance will not increase under these maps. Now describing the states with the high Hamming weight components removed by $|\tilde{\psi}_\eta\rangle$, we have

$$\begin{aligned} & \left\| \sum_{|\mathbf{b}| \leq k'} \left[|\mathbf{b}\rangle\langle\mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\eta, \mathbf{b}} |\tilde{\psi}_\eta\rangle\langle\tilde{\psi}_\eta| M_{\eta, \mathbf{b}}^\dagger \right) - |\mathbf{b}\rangle\langle\mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\eta, \mathbf{b}} |\psi_\eta\rangle\langle\psi_\eta| M_{\eta, \mathbf{b}}^\dagger \right) \right] \right\|_{\text{tr}} \\ & \leq \left\| \mathcal{E}_\eta(|\tilde{\psi}_\eta\rangle\langle\tilde{\psi}_\eta|) - \mathcal{E}_\eta(|\psi_\eta\rangle\langle\psi_\eta|) \right\|_{\text{tr}} \leq \left\| |\tilde{\psi}_\eta\rangle\langle\tilde{\psi}_\eta| - |\psi_\eta\rangle\langle\psi_\eta| \right\|_{\text{tr}} = O(\varepsilon). \end{aligned} \quad (\text{B.8})$$

As a result, using the triangle inequality gives

$$\begin{aligned} D'_{\text{av}} & \leq O(\varepsilon) \\ & + \left\| \sum_{|\mathbf{b}| \leq k'} \left[|\mathbf{b}\rangle\langle\mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\mathbf{u}, \mathbf{b}} |\tilde{\psi}_{\mathbf{u}}\rangle\langle\tilde{\psi}_{\mathbf{u}}| M_{\mathbf{u}, \mathbf{b}}^\dagger \right) - |\mathbf{b}\rangle\langle\mathbf{b}| \otimes \text{Tr}_{\text{ctrl}} \left(M_{\mathbf{c}, \mathbf{b}} |\tilde{\psi}_{\mathbf{c}}\rangle\langle\tilde{\psi}_{\mathbf{c}}| M_{\mathbf{c}, \mathbf{b}}^\dagger \right) \right] \right\|_{\text{tr}}. \end{aligned} \quad (\text{B.9})$$

Next, each measurement operator $M_{\eta, \mathbf{b}}$ can be obtained by a sequence of measurement operators in our recursive measurement scheme, which will yield a sequence of measurement results d_1, d_2, \dots . Each \mathbf{b} will correspond to a unique sequence of d_ℓ measurement results. (Recall that b_j are the individual results of measurements on uncompressed qubits, whereas d_ℓ are the individual results from the recursive measurement.) Therefore we can relabel the basis states such that we have

$$\begin{aligned} D'_{\text{av}} & = \left\| \sum_{\mathbf{d}} \left\{ |\mathbf{d}\rangle\langle\mathbf{d}| \otimes \text{Tr}_{\text{ctrl}} \left[M'_{\mathbf{u}, \mathbf{d}} |\psi_{\mathbf{u}}\rangle\langle\psi_{\mathbf{u}}| (M'_{\mathbf{u}, \mathbf{d}})^\dagger \right] - |\mathbf{d}\rangle\langle\mathbf{d}| \otimes \text{Tr}_{\text{ctrl}} \left[M'_{\mathbf{c}, \mathbf{d}} |\psi_{\mathbf{c}}\rangle\langle\psi_{\mathbf{c}}| (M'_{\mathbf{c}, \mathbf{d}})^\dagger \right] \right\} \right\|_{\text{tr}}. \end{aligned} \quad (\text{B.10})$$

Now the measurement operators that are chosen at step ℓ in the recursive measurement scheme will depend on the measurement results that have been obtained at steps 1 to $\ell - 1$. Therefore we can write the measurement operators as

$$M'_{\eta, \mathbf{d}} = \prod_{\ell=1}^K M_{\eta, d_\ell}^{(d_1, \dots, d_{\ell-1})}. \quad (\text{B.11})$$

Here K is the number of measurement operators to locate the ones. For measurement result \mathbf{b} , the number of measurements required is no more than $1 + 2|\mathbf{b}| \log m$. As we are taking \mathbf{b} such that $|\mathbf{b}| \leq k'$, we can take $K = 1 + 2k' \log m$.

Using this notation, we can define CPTP maps by

$$\mathcal{E}_{\eta, \ell}(\rho) := \sum_{d_1, \dots, d_\ell} |d_\ell\rangle\langle d_\ell| \otimes M_{\eta, d_\ell, \text{proj}}^{(d_1, \dots, d_{\ell-1})} \rho (M_{\eta, d_\ell, \text{proj}}^{(d_1, \dots, d_{\ell-1})})^\dagger, \quad (\text{B.12})$$

where

$$M_{\eta, d_\ell, \text{proj}}^{(d_1, \dots, d_{\ell-1})} := |d_{\ell-1}\rangle\langle d_{\ell-1}| \otimes \dots \otimes |d_1\rangle\langle d_1| \otimes M_{\eta, d_\ell}^{(d_1, \dots, d_{\ell-1})}. \quad (\text{B.13})$$

Each map simply performs the appropriate measurement based on the prior measurement results (which are stored in ancillas), and appends an ancilla depending on the result of the measurement. In term of these maps, the trace distance we wish to bound may be written as

$$D'_{\text{av}} = \|\text{Tr}_{\text{ctrl}} \mathcal{E}_{\text{u},K} \dots \mathcal{E}_{\text{u},1}(|\psi_{\text{u}}\rangle\langle\psi_{\text{u}}|) - \text{Tr}_{\text{ctrl}} \mathcal{E}_{\text{c},K} \dots \mathcal{E}_{\text{c},1}(|\psi_{\text{c}}\rangle\langle\psi_{\text{c}}|)\|_{\text{tr}}. \quad (\text{B.14})$$

As has been noted above, we can omit the high Hamming weight contributions to the states $|\psi_{\eta}\rangle$, with a possible change in the trace distance of $O(\varepsilon)$. The reason for this is that the trace distance is non-increasing under CPTP maps. Our goal is now to successively approximate each of the maps in the sequence, at each stage bounding the introduced error by $O(\varepsilon)$. At the end we will obtain two identical states, and then bound D'_{av} by $O(K\varepsilon)$.

More specifically, we want to approximate the evolution of the states for given measurement results as in Eqs. (38) and (39). Note that the reasoning given in the proof of Theorem 5 also gives a recursive method to determine the amplitudes in the states $|\psi_{\eta,\ell}^{(d_1,\dots,d_{\ell})}\rangle$, starting from $\gamma_{x,0}^{(0)} = i^{|x|}\xi_x^m$. This means that the definitions of these states are unambiguous. We now consider the approximate unnormalised states after $\ell - 1$ measurements in the recursive measurement scheme $|\psi_{\eta,\ell-1}^{(d_1,\dots,d_{\ell-1})}\rangle$, as given by Eqs. (35) and (36). We then define the states including the ancilla qubits containing the measurement results as

$$\rho_{\eta,\ell-1} := \sum_{d_1,\dots,d_{\ell-1}} |d_{\ell-1}\rangle\langle d_{\ell-1}| \otimes \dots \otimes |d_1\rangle\langle d_1| \otimes |\psi_{\eta,\ell-1}^{(d_1,\dots,d_{\ell-1})}\rangle\langle\psi_{\eta,\ell-1}^{(d_1,\dots,d_{\ell-1})}|. \quad (\text{B.15})$$

We wish to bound the error in approximating $\mathcal{E}_{\eta,\ell}(\rho_{\eta,\ell-1})$ by $\rho_{\eta,\ell}$.

In approximating $\mathcal{E}_{\text{u},\ell}(\rho_{\text{u},\ell-1})$ by $\rho_{\text{u},\ell}$ there is only one approximation: that of omitting the high Hamming weight states in applying the rotation. The error in this approximation will be $O(\varepsilon)$ times the norm of the state. Because the norm of the state is only changed by omitting high Hamming weight components, it can only be decreased. Therefore the error is $O(\varepsilon)$. Similarly, there is error in approximating $\mathcal{E}_{\text{c},\ell}(\rho_{\text{c},\ell-1})$ by $\rho_{\text{c},\ell}$ due to omitting high Hamming weight components, which is bounded by $O(\varepsilon)$. There is also error because the U_n rotations are not performed exactly. Two such rotations are performed, each with error bounded by $O(\varepsilon)$, resulting in the overall error being bounded by $O(\varepsilon)$.

Therefore, we can start with Eq. (B.10), remove the high Hamming weight components from the initial states, then proceed taking $\ell = 1$ to K , replacing $\mathcal{E}_{\eta,\ell}(\rho_{\eta,\ell-1})$ by $\rho_{\eta,\ell}$ at each step. At each step the distance is increased by $O(\varepsilon)$, and there are K steps, so we obtain

$$D'_{\text{av}} \leq O(K\varepsilon) + \|\text{Tr}_{\text{ctrl}}(\rho_{\text{u},K}) - \text{Tr}_{\text{ctrl}}(\rho_{\text{c},K})\|_{\text{tr}}. \quad (\text{B.16})$$

But, because the same amplitudes have been obtained for the compressed and uncompressed cases, the same state is obtained after tracing over the control registers, and $\text{Tr}_{\text{ctrl}}(\rho_{\text{u},K}) = \text{Tr}_{\text{ctrl}}(\rho_{\text{c},K})$. Therefore we obtain

$$D'_{\text{av}} = O(K\varepsilon) = O(\varepsilon k' \log m). \quad (\text{B.17})$$

As $D_{\text{av}} = O(\varepsilon' + D'_{\text{av}})$, this yields Eq. (45), as required.