

# Memory-Sample Lower Bounds for Learning with Classical-Quantum Hybrid Memory

Qipeng Liu <sup>\*</sup>      Ran Raz <sup>†</sup>      Wei Zhan <sup>‡</sup>

## Abstract

In a work by Raz (J. ACM and FOCS 16), it was proved that any algorithm for parity learning on  $n$  bits requires either  $\Omega(n^2)$  bits of classical memory or an exponential number (in  $n$ ) of random samples. A line of recent works continued that research direction and showed that for a large collection of classical learning tasks, either super-linear classical memory size or super-polynomially many samples are needed. All these works consider learning algorithms as classical branching programs, which perform classical computation within bounded memory.

However, these results do not capture all physical computational models, remarkably, quantum computers and the use of quantum memory. It leaves the possibility that a small piece of quantum memory could significantly reduce the need for classical memory or samples and thus completely change the nature of the classical learning task. Despite the recent research on the necessity of quantum memory for intrinsic quantum learning problems like shadow tomography and purity testing, the role of quantum memory in classical learning tasks remains obscure.

In this work, we study classical learning tasks in the presence of quantum memory. We prove that any quantum algorithm with both, classical memory and quantum memory, for parity learning on  $n$  bits, requires either  $\Omega(n^2)$  bits of classical memory or  $\Omega(n)$  bits of quantum memory or an exponential number of samples. In other words, the memory-sample lower bound for parity learning remains qualitatively the same, even if the learning algorithm can use, in addition to the classical memory, a quantum memory of size  $cn$  (for some constant  $c > 0$ ).

Our result is more general and applies to many other classical learning tasks. Following previous works, we represent by the matrix  $M : A \times X \rightarrow \{-1, 1\}$  the following learning task. An unknown  $x$  is sampled uniformly at random from a concept class  $X$ , and a learning algorithm tries to uncover  $x$  by seeing streaming of random samples  $(a_i, b_i = M(a_i, x))$  where for every  $i$ ,  $a_i \in A$  is chosen uniformly at random. Assume that  $k, \ell, r$  are integers such that any submatrix of  $M$  of at least  $2^{-k} \cdot |A|$  rows and at least  $2^{-\ell} \cdot |X|$  columns, has a bias of at most  $2^{-r}$ . We prove that any algorithm with classical and quantum hybrid memory for the learning problem corresponding to  $M$  needs either (1)  $\Omega(k \cdot \ell)$  bits of classical memory, or (2)  $\Omega(r)$  qubits of quantum memory, or (3)  $2^{\Omega(r)}$  random samples, to achieve a success probability at least  $2^{-O(r)}$ .

Our results refute the possibility that a small amount of quantum memory significantly reduces the size of classical memory needed for efficient learning on these problems. Our results also imply improved security of several existing cryptographical protocols in the bounded-storage model (protocols that are based on parity learning on  $n$  bits), proving that security holds even in the presence of a quantum adversary with at most  $cn^2$  bits of classical memory and  $cn$  bits of quantum memory (for some constant  $c > 0$ ).

---

<sup>\*</sup>Simons Institute for Theory of Computing. E-mail: qipengliu0@gmail.com. Research supported in part by the Simons Institute for the Theory of Computing, through a Quantum Postdoctoral Fellowship, by the DARPA SIEVE-VESPA grant No.HR00112020023 and by the NSF QLCI program through grant number OMA-2016245.

<sup>†</sup>Princeton University. E-mail: ranr@cs.princeton.edu. Research supported by a Simons Investigator Award and by the National Science Foundation grants No. CCF-1714779, CCF-2007462.

<sup>‡</sup>Princeton University. E-mail: weizhan@cs.princeton.edu. Research supported by a Simons Investigator Award and by the National Science Foundation grants No. CCF-1714779, CCF-2007462.

# 1 Introduction

Memory plays an important role in learning. Starting from the seminal works by Shamir [Sha14] and Steinhardt, Valiant and Wager [SVW16], a sequence of works initiates and deepens the study of lower bounds for learning under memory constraints. Steinhardt, Valiant, and Wager [SVW16] conjectured that in order to learn an unknown  $n$ -bit string from samples of random-subset parity, an algorithm needs either memory-size quadratic in  $n$  or exponentially many random samples (also in  $n$ ). This conjecture was later on proved by Raz [Raz18], showing for the first time that for some learning problems, super-linear memory size is required for efficient learning. This result was then generalized to a broad class of learning problems [KRT17, Raz17, MM18, BGY18, GRT18, SSV19, GKR20, GKLR21].

Although we have a comprehensive understanding of the (in)feasibility of learning under limitations on particular computation resource (memory), the previous works mentioned above do not capture all physical computational models; most notably, quantum computation and the power of quantum memory. Many researchers believe that large-scale quantum computers will eventually become viable. Recent experiments demonstrated quantum advantages, for example [AAB<sup>+</sup>19], and suggested that there are possibly no fundamental barriers to achieving quantum memory and quantum computers. Questions on the role of quantum memory in learning were proposed by Wright in the context of general state tomography [Wri16] and by Aaronson for shadow tomography [Aar18]. A line of works [HHJ<sup>+</sup>16, BCL20, HKP20, CCHL22, ACQ22, CLO22] pioneer the idea and show either polynomial or exponential separations for learning with/without quantum memory, but all for intrinsic quantum learning tasks like state tomography, shadow tomography and purity testing.

In light of the above, it is appealing to consider classical learning tasks in the presence of quantum memory, as well as hybrid classical-quantum memory. A direct implication of all aforementioned classical results only gives trivial results. As  $k$  qubits of memory can always be efficiently simulated by  $\sim 2^k$  classical bits, we can only conclude (say, for parity learning) that either  $\sim 2 \log n$ -qubit quantum memory or exponentially many samples are needed. Prior to our work, it could have been the case that even if only a very small size quantum memory was available, it might have significantly reduced the need for classical memory and led to an efficient learning algorithm.

In this work, we prove memory-sample lower bounds in the presence of hybrid memory for a wide collection of classical learning problems. As in [Raz17, GRT18], we will represent a learning problem by a matrix  $M : A \times X \rightarrow \{-1, 1\}$  whose columns correspond to concepts in the concept class  $X$  and rows correspond to random samples. In the learning task, an unknown concept  $x \in X$  is sampled uniformly at random and each random sample is given as  $(a_i, b_i) = (a_i, M(a_i, x))$  for a uniformly picked  $a_i \in A$ . The learner's goal is to uncover  $x$ . In [GRT18], it is proved that when the underlying matrix  $M$  is a  $(k, \ell)$ - $L_2$  two source extractor<sup>1</sup> with error  $2^{-r}$ , a learning algorithm requires either  $\Omega(k \cdot \ell)$  bits of memory or  $2^{\Omega(r)}$  samples to achieve a success probability at least  $2^{-O(r)}$  for the learning task.

## 1.1 Our Results

In this work, we model a quantum learning algorithm as a program with hybrid memory consisting of  $q$  qubits of quantum memory and  $m$  bits of classical memory. At each stage, a random sample  $(a_i, b_i = M(a_i, x))$  is given to the algorithm. The quantum learning algorithm applies an arbitrary quantum channel to the hybrid memory, controlled by the random sample. Although the channel

---

<sup>1</sup>Roughly speaking, this means that every submatrix  $M'$  of  $M$  with number of rows at least  $2^{-k}|A|$  and number of columns at least  $2^{-\ell}|X|$  has a relative bias at most  $2^{-r}$ .

can be arbitrary, we impose the outcome to be a hybrid classical-quantum state of at most  $q$  qubits and  $m$  bits. We stress that there is no limitation on the complexity of the quantum channel (and this only makes our results stronger as we are proving here lower bounds for such algorithms).

With the above model, we give the following main theorem.

**Theorem 1** (Main Theorem, Informal). *Let  $M : A \times X \rightarrow \{-1, 1\}$  be a matrix. If  $M$  is a  $(k, \ell)$ - $L_2$  two source extractor with error  $2^{-r}$ , a quantum learning algorithm requires either*

1.  $\Omega(k \cdot \ell)$  bits of classical memory; or,
2.  $\Omega(r)$  qubits of quantum memory; or,
3.  $2^{\Omega(r)}$  samples,

*to succeed with a probability of at least  $2^{-O(r)}$  in the corresponding learning task.*

Our main theorem implies that for many learning problems, the availability of a quantum memory of size up to  $\Omega(r)$ , does not reduce the size of classical memory or the number of samples that are needed. As coherent quantum memory is challenging for near-term intermediate-scale quantum computers and is probably expensive even if and when quantum computers are widely viable, the impact of quantum memory is further limited for these learning problems.

To make the theorem more precise, let us take parity learning as an example. The above theorem says that a quantum learning algorithm needs either  $\Omega(n^2)$  bits of memory, or  $\Omega(n)$  qubits of quantum memory, to conduct efficient learning; otherwise, it requires  $2^{\Omega(n)}$  random samples. At first glance, it seems that the constraint on quantum memory is trivial: if the target is to learn an  $n$ -bit unknown secret, a linear amount of memory always seems necessary to store the secret. However, noticing that our main theorem applies to quantum learning algorithms with hybrid memory and rules out algorithms with  $n^2/1000$  bits and  $n/1000$  qubits of hybrid memory for parity learning, the main theorem yields non-trivial and compelling memory-sample lower bounds. Note also that our results (and previous results) are valid even if the goal is to output only one bit of the secret. Currently, we do not know whether our main theorem is tight. For parity learning, we are not aware of any quantum learning algorithm that uses only  $O(n)$  qubits of quantum memory. We leave closing the gap as a fascinating open question.

The main theorem naturally applies to other learning problems considered in [GRT18], including learning sparse parities, learning from sparse linear equations, and many others. We do not present an exhaustive list here but refer the readers to [GRT18] for more details.

Along the way, we propose a new approach for proving the classical memory-sample lower bounds. We call this approach, the “badness levels” method. The approach is technically equivalent to the previous approach in [Raz17, GRT18] but is conceptually simpler to work with and we are able to lift it to the quantum case.

We note that proving a linear lower bound on the size of the quantum memory, without classical memory, is significantly simpler (but to the best of our knowledge such a proof has not appeared prior to our work). We present such a proof in Appendix C. In Appendix C, we state and prove Theorem 3 that shows a simpler proof for a linear lower bound on the quantum-memory size (without classical memory). While Theorem 3 is qualitatively weaker than our main result in most cases, as it only gives a lower bound for programs with only quantum memory but without a (possibly quadratic) classical memory, Theorem 3 is technically incomparable and is stated in terms of quantum extractors, rather than classical extractors. Additionally, the proof of Theorem 3 is significantly simpler than the proof of our main theorem.

**Implications to Cryptography in the Bounded-Storage Model.** Since learning theory and cryptography can be viewed as two sides of the same coin, our theorem also lifts the security of many existing cryptographical protocols in the bounded-storage model (protocols that are based on parity learning) to the quantum setting. To our best knowledge, these are the first proofs of classical cryptographical protocols being secure against space-bounded quantum attackers.<sup>2</sup> We elaborate more below.

Cryptography in the (classical) bounded storage model was first proposed by Maurer [Mau92]. In such a model, no computational assumption is needed. Honest execution is performed through a long streaming sequence of bits. Eavesdroppers have bounded storage and limited capability of storing conversations, thus cannot break the protocol. A line of works [CM97, AR99, ADR02, DR02, DM02, Lu02, DM04, MST04, HN06, DQW21, DQW22, ...] builds efficient and secure protocols for key agreement, oblivious transfer, bit commitment and time stamping in that model.

Based on the memory-sample lower bounds for parity learning of  $n$  bits, [Raz18] suggested an encryption scheme in the bounded-storage model. Guan and Zhandry [GZ19] proposed key agreement, oblivious transfer and bit commitment with improved rounds and better correctness, against attackers with up to  $O(n^2)$  bits of memory. Following a similar idea, Liu and Vusirikala [LV21] showed that semi-honest multiparty computation could be achieved against attackers with up to  $O(n^2)$  bits of memory. More recently, Dodis, Quach, and Wichs [DQW22] considered message authentication in the bounded storage model based on parity learning. Our result on parity learning gives a direct lift on all the results above. When the cryptographic protocols are based on parity learning of  $n$  bits (often treated as a security parameter), our result shows that security holds even in the presence of a quantum adversary with at most  $O(n^2)$  bits of classical memory and  $O(n)$  qubits of quantum memory.

Despite many previous works on cryptography in the quantum bounded storage model [DFR<sup>+</sup>07, DFSS07, Sch07, DFSS08, WW08, PMLA13, BY21], they all rely on streaming quantum states. Our memory-sample lower bounds give for the first time a rich class of classical cryptographical schemes (key agreement, oblivious transfer, and bit commitment) secure against space-bounded quantum attackers.

## 2 Proof Overview

### 2.1 Recap of Proofs for Classical Lower Bounds

Since our proof builds on the previous line of works on classical memory-sample lower bounds for learning, specifically, on the proof technique of [Raz17, GRT18], we provide a brief review of these proofs, using parity learning [Raz18] as an example. In below,  $M(a, x)$  denotes the inner product of  $a$  and  $x$  in  $\mathbb{F}_2$ .

Consider a classical branching program that tries to learn an unknown and uniformly random  $x \in \{0, 1\}^n$  from samples  $(a, b)$ , where  $a \in \{0, 1\}^n$  is uniformly random and  $b = M(a, x)$ . We can associate every state  $v$  of the branching program with a distribution  $P_{X|v}$  over  $\{0, 1\}^n$ , indicating the distribution of  $x$  conditioned on reaching that state. At the initial state, without any information about  $x$ , the distribution is uniform (which has the smallest possible  $\ell_2$ -norm). Along a computational path on the branching program, the distribution  $P_{X|v}$  evolves and eventually gets concentrated (with large  $\ell_2$ -norms) in order to output  $x$  correctly. Therefore, during the evolution,  $P_{X|v}$  should at some stage have mildly large  $\ell_2$ -norms ( $2^{\epsilon n}$  times larger than uniform for some

---

<sup>2</sup>On the other hand, there are known examples of classically-secure bounded-storage protocols that are breakable with an exponentially smaller amount of quantum memory. [GKK<sup>+</sup>08].

small constant  $\varepsilon > 0$ ). If we set such a distribution as a target, the distribution is hard to achieve with random samples. Only with  $2^{-\Omega(n)}$  probability, the branching program can make significant progress towards the target; while most of the time a sample just splits the distributions (both the current and the target distribution) into two even parts, and that does not help much in getting closer to the target distribution (with large  $\ell_2$  norm).

To put it more rigorously, we examine the evolution of the inner product

$$\langle P_{X|v}, P \rangle = \sum_{x \in \{0,1\}^n} P_{X|v}(x) \cdot P(x)$$

between the distribution  $P_{X|v}$  on the current state  $v$ , and a target distribution  $P$ . Receiving a sample  $(a, b)$  implies that  $M(a, x) = b$ , hence only the part of  $P_{X|v}$  supported on such  $x$  proceeds. If this part is close to  $\frac{1}{2}$  probability, we say that  $a$  divides  $P_{X|v}$  evenly. Denoting the new distribution as  $P_{X|v}^{(a,b)}$ , after proper normalization the new inner product is

$$\langle P_{X|v}^{(a,b)}, P \rangle = \sum_{\substack{x \in \{0,1\}^n \\ M(a,x)=b}} P_{X|v}(x) \cdot P(x) \Big/ \sum_{\substack{x \in \{0,1\}^n \\ M(a,x)=b}} P_{X|v}(x). \quad (1)$$

Ideally, both  $P_{X|v}$  and the point-wise product vector  $P_{X|v} \cdot P$  should have reasonably small  $\ell_2$ -norms. Due to the extractor property of  $M$ , most of  $a \in \{0,1\}^n$  should divide both vectors evenly, and thus the denominator is close to  $\frac{1}{2}$  while the numerator is close to  $\frac{1}{2} \langle P_{X|v}, P \rangle$ . That means, given a uniformly random  $a$ , we get limited progress on the inner product. On the other hand, from  $\langle U, P \rangle = 2^{-n}$  with uniform distribution  $U$  to  $\langle P, P \rangle = 2^{2\varepsilon n} \cdot 2^{-n}$ , the branching program needs to make multiple steps of progression. Therefore it happens with an extremely small probability.

To ensure that the above statement goes smoothly, we require the following properties for every state  $v$  in the branching program:

- The  $\ell_2$ -norm  $\|P_{X|v}\|_2$  is small.
- The  $\ell_2$ -norm  $\|P_{X|v} \cdot P\|_2$  is small, which is implied when the  $\ell_\infty$ -norm  $\|P_{X|v}\|_\infty$  is small.
- The denominator in Eq. (1) is bounded away from 0 for every sample  $(a, b)$ .

These properties do not hold by themselves. Instead, we execute a *truncation* procedure on the branching program *before* choosing a target distribution. More specifically, the branching program is modified so that it stops whenever it:

- ( $\ell_2$  truncation): Reaches a state  $v$  with large  $\|P_{X|v}\|_2$ ;
- ( $\ell_\infty$  truncation): Reaches a state  $v$  with large  $P_{X|v}(x)$  when the unknown concept is  $x$ ;
- (Sample truncation): Or, for the next sample  $(a, b)$ ,  $a$  does not divide  $P_{X|v}$  evenly.

It turns out that after  $\ell_2$  truncation, the other two truncation steps add  $2^{-\Omega(n)}$  error in each stage of the branching program. Therefore the proof boils down to proving a  $2^{-\Omega(n^2)}$  bound on the probability of reaching a state with large  $\|P_{X|v}\|_2$ , from which by a standard union bound, we can prove the memory-sample lower bounds for parity learning: either  $2^{\Omega(n)}$  samples or  $\Omega(n^2)$  bits of memory are necessary.

## 2.2 Badness Levels

As mentioned above, to bound the probability of reaching a state with a large  $\ell_2$ -norm, the basic idea is to fix its distribution as the target distribution  $P$ , and bound the increment of the inner product  $\langle P_{X|v}, P \rangle$ . This was done in [Raz17, GRT18] by designing a potential function that tracks the average of  $\langle P_{X|v}, P \rangle^k$  for some  $k = \Theta(n)$ , where the average is over states  $v$  in the same stage of the branching program. Here we propose another approach using the concept of *badness levels*. Although it is technically equivalent to the potential function approach in the classical case, it is more pliable and easier to be adapted to the quantum case. We view this approach as a separate contribution of our work.

We first define a *bad event* to be a pair  $(v, a)$  of the state  $v$  and the upcoming part of the sample  $a$ , such that  $\langle P_{X|v}, P \rangle \geq 2^{-n}$ , and for one of the two possible outcomes  $b$ ,

$$\sum_{\substack{x \in \{0,1\}^n \\ M(a,x)=b}} P_{X|v}(x) \cdot P(x) \geq \left(\frac{1}{2} + 2^{-\delta n}\right) \cdot \langle P_{X|v}, P \rangle \quad (2)$$

with some small constant  $\delta$ . In other words, the inner product  $\langle P_{X|v}, P \rangle$  is large enough, while not being divided evenly by  $a$ . From Eq. (1) we know that the inner product gets at most roughly doubled through a bad event. In contrast, in the good case, the inner product either gets a mere  $(1 + 2^{-\delta n})$  multiplicative factor or is already smaller than the baseline  $2^{-n}$ . Also, the extractor property of  $M$  ensures that for every state  $v$ , over uniformly random  $a$ , the bad event happens with at most  $2^{-\Omega(n)}$  probability.

Now, the badness level  $\beta(v)$  of a state  $v$  keeps track of how many times the computational path went through bad events before reaching  $v$ .<sup>3</sup> The above observations on the bad events imply that (omitting the smaller factors):

- For every state  $v$ ,  $\langle P_{X|v}, P \rangle$  is bounded by  $2^{\beta(v)} \cdot 2^{-n}$ ;
- Heading to the next stage,  $\beta(v)$  increases by 1 with probability  $2^{-\Omega(n)}$ .

Therefore at each stage, the total weight of states with badness level  $\beta$  is at most  $2^{-\Omega(\beta n)}$ . Thus any state with  $\langle P_{X|v}, P \rangle \geq 2^{2\epsilon n} \cdot 2^{-n}$  must have  $2^{-\Omega(n^2)}$  probability.

## 2.3 Obstacles for Proving Quantum Lower Bounds

In this section, we present an attempt to prove the same  $2^{\Omega(n)}$ -sample or  $\Omega(n^2)$ -quantum-memory lower bound for the pure quantum case. Along the way we identify some obstacles to proving memory-sample lower bounds for *quantum* learning algorithms, and in the next section we show how to overcome these obstacles while proving lower bounds for *hybrid* learning algorithms, with quadratic-size classical-memory and linear-size quantum-memory.

Following the same framework as the above described proof for the classical case, we first need to transfer all the notions to a quantum algorithms:

- The state  $v$  is a quantum state in the Hilbert space of quantum memory;
- The distribution  $P_{X|v}$  is still well-defined: It is the distribution of  $x$  when the quantum memory is measured to  $v$  (see Section 3.4 and Eq. (3));

---

<sup>3</sup>For now we think of  $\beta(v)$  as a natural number. In the actual proof,  $\beta(v)$  is a distribution on natural numbers, as for different computational paths reaching the same state, the count of bad events can be different.



- We are still able to implement  $\ell_2$  truncation: If  $P_{X|v}$  has large  $\ell_2$ -norm, project the entire system to the orthogonal subspace  $v^\perp$  of  $v$  and repeat, until there is no such state  $v$  (see Section 4.1 for details).
- We are also able to implement sample truncation, in a similar manner to  $\ell_2$  truncation. As the criteria here depends on  $a$ , we separately create a copy of the current system for each  $a$ , truncate the states  $v$  using projection when  $P_{X|v}$  is not evenly divided by  $a$  in each copy, and then merge them back together. We prove that the error introduced by this truncation is small.

Here comes the first major obstacle:  $\ell_\infty$  truncation. In the classical case,  $\ell_\infty$  truncation is implemented for each individual  $x$ , in contrast to  $\ell_2$  truncation where the states are removed altogether. Relying on the fact that it is already known that the  $\ell_2$  norm of the distribution is small, using Markov inequality, one can prove that the error introduced by the  $\ell_\infty$  truncation is small.

However, when we try to emulate the classical implementation of  $\ell_\infty$  truncation with quantum truncation, that is, to only project to  $v^\perp$  the system *conditioned on* the specific  $x$  where  $P_{X|v}(x)$  is large, instead of for every  $x$ , it may lead to huge changes to the distributions  $P_{X|u}$  on states  $u$  non-orthogonal to  $v$ . The following example illustrates such a scenario:

**Example.** Consider a quantum learning algorithm, and assume that at some stage of the computation, for each  $x \in \{0, 1\}^n$ , the quantum memory is in some pure state  $v(x)$ . We pick each  $v(x)$  uniformly at random in a Hilbert space of dimension  $d \approx 2^{n/2}$  and consider a typical configuration of  $v(x)$ . Now the  $\ell_2$ -norms are bounded for every quantum state  $v$ : the worst ones happen when  $v = v(x)$  for some  $x$ , where  $\|P_{X|v(x)}\|_2$  is typically around  $d \cdot 2^{-n}$ , close to the  $\ell_2$ -norm of uniform distribution. However, those worst distributions also have  $\ell_\infty$ -norms close to  $d \cdot 2^{-n}$ , which is much larger than the  $\ell_\infty$ -norm of the uniform distribution, and needs to be truncated. But truncating  $v(x)$  off for  $x$  means that  $x$  is completely erased, and we end up removing everything.

Moving on, we fix a target state  $v$  with a target distribution  $P_{X|v}$  which exceeds the  $\ell_2$ -norm threshold, and the goal is again to prove a  $2^{-\Omega(n^2)}$  amplitude bound on  $v$ . The bad event should still be defined as a pair  $(v, a)$  satisfying Eq. (2), with  $v$  now being a quantum state. We then run into the second major obstacle: it is not clear how to define badness levels.

If we define the badness level  $\beta(v)$  for each state  $v$  individually by examining the bad events over the historical states, then it is not clear how to measure the total weight of a badness level  $\beta$ . In the classical case, we simply define the total weight as the total probability of states with badness level  $\beta$ . But here in the quantum case, it turns out that such a definition either depends on the choice of basis, which might have large increment in each stage, or completely fails to imply the desired amplitude bound on the target state.

The other choice is to have a more *operational* definition of badness levels, and it is indeed tempting to define  $\beta$  as another register whose updates are controlled by the quantum memory. The problem with such definitions is that the bad event (Eq. (2)) is not linear in  $v$ . Therefore an operational definition of badness level, which is a linear operator, inevitably introduces error that escalates fast with the number of stages.

## 2.4 Hybrid Memory Lower Bounds with Small Quantum Memory

The obstacles in the previous section are for proving quadratic quantum memory lower bound. We note that proving linear quantum memory lower bound (without classical memory) is not hard:

the proof can be entirely information theoretical, as with very limited memory, say,  $\frac{1}{2}n$  qubits, the information gained from each sample is exponentially small, despite the memory being quantum. We present such a proof in Appendix C.

The lower bounds that we prove here are with hybrid memory: To learn parity with both classical and quantum memory, an algorithm needs either  $2^{\Omega(n)}$  samples, or  $\Omega(n^2)$  classical memory, or  $\Omega(n)$  quantum memory (Theorem 1). We now describe how we overcome the previously mentioned obstacles.

**$\ell_\infty$  Truncation.** When there is only small quantum memory and no classical memory, the treatment for  $\ell_\infty$  truncation is straightforward. We remove all quantum states  $v$  with distributions of large  $\ell_\infty$ -norm, by projecting the system to the orthogonal subspace  $v^\perp$ , just like the process of  $\ell_2$  truncation. As the overall distribution on  $x$  is uniform, any state  $v$  with  $\|P_{X|v}\|_\infty \geq 2^{\delta n} \cdot 2^{-n}$  must have weight at most  $2^{-\delta n}$ . Therefore, as long as the dimension of the Hilbert space is much smaller than  $\delta n$ , the error introduced in this truncation is small.<sup>4</sup>

With classical memory in presence, the actual  $\ell_\infty$  truncation step (see Section 4.2, Step 2) is more complicated. We first apply the original classical  $\ell_\infty$  truncation on the classical memory  $W$ . Now that  $\|P_{X|w}\|_\infty$  is bounded for each classical memory state  $w$ , we can remove the quantum states  $v$  with large  $\|P_{X|v,w}\|_\infty$  by projection as stated above. Since the classical  $\ell_\infty$  truncation depends on  $x$ , it could change the distributions  $P_{X|v,w}$ . However, as in the classical case,  $P_{X|w}$  will not change a lot. Thus, wherever  $P_{X|v,w}$  changes drastically, it must have a small weight and can also be removed by projection. This removal corresponds to truncation by  $G_t$  in Section 4.2.

**Badness Levels.** Interestingly, we are able to avoid the problems of defining the badness level on quantum memory altogether, by keeping it a property on the classical memory only. To do so we need to alter the definition of a bad event: it is now a pair  $(w, a)$  of classical memory state  $w$  and sample  $a$ , such that *there exists* some quantum memory state  $v$  with  $P_{X|v,w}$  satisfying Eq. (2).

For each fixed classical memory state  $w$ , we still need to ensure that bad events happen with a small probability. We prove it (Lemma 5.2) by showing that, if there are many different samples  $a$ , each associated with some quantum state  $v_a$  satisfying Eq. (2), then there is some quantum state  $v$  that simultaneously satisfies Eq. (2) with most of such  $a$  (which is impossible because of the extractor property). This is ultimately due to the continuous nature of Eq. (2): Under some proper congruent transformation, Eq. (2) becomes a simple threshold inequality on quadratic forms over  $v$ . Now if it is satisfied by some  $v_a$ , it is going to be satisfied by most  $v$  for a much smaller threshold parameter  $\delta$ , and hence the existence of a simultaneously satisfying  $v$ .<sup>5</sup> In this argument, we use Lemma 3.1, which is derived from the anti-concentration bound for Gaussian quadratic forms, and crucially relies on the fact that the dimension is at most  $2^{\epsilon n}$  for some small  $\epsilon$ .

Another technical problem is that to use the extractor property, we need to ensure that  $\langle P_{X|v,w}, P \rangle \geq 2^{-n}$  for the simultaneously satisfying  $v$ . Thus, what we do in Lemma 5.2 is to first conceptually remove the parts where  $\langle P_{X|v,w}, P \rangle$  is too small, using projection similarly to the truncation steps. After the removal, we are left with a subspace  $\mathcal{V}'$  where  $\langle P_{X|v,w}, P \rangle$  is always lower bounded, and we show that for every state  $v$  that satisfies Eq. (2), the inequality is still close to being satisfied after projecting  $v$  onto  $\mathcal{V}'$ . Therefore we could still apply the above argument and find a simultaneously satisfying  $v$  within the subspace.

<sup>4</sup>The example in the previous section that shows the infeasibility of treating  $\ell_\infty$  truncation the same way as  $\ell_2$  truncation does not work here, as it requires  $n/2$  qubits of memory while here we have a smaller memory size.

<sup>5</sup>We note that the error bound for sample truncation (Lemma 4.12) is also proved using this argument.



### 3 Preliminaries

#### 3.1 Vectors and Matrices

For a vector  $v \in \mathbb{C}^d$  and  $p \in [1, \infty]$ , we define the  $\ell_p$  norm of  $v$  as

$$\|v\|_p = \left( \sum_{i=1}^d |v_i|^p \right)^{1/p}.$$

For two vectors  $u, v \in \mathbb{C}^d$ , define their inner product as  $\langle u, v \rangle = u^\dagger v = \sum_{i=1}^d \overline{u_i} v_i$ . So  $\|v\|_2^2 = \langle v, v \rangle$ . We also view every distribution  $P$  over a set  $\mathcal{X}$  as a non-negative real vector with  $\|P\|_1 = 1$ .

We specifically use Dirac notation to denote unit vectors,  $|v\rangle \in \mathbb{C}^d$  implies that  $\| |v\rangle \|_2 = 1$ . For a non-zero vector  $u \in \mathbb{C}^d$ , let  $|v\rangle \sim u$  be the normalization of  $u$ , that is,  $|v\rangle = u/\|u\|_2$ .

For every vector  $v \in \mathbb{C}^d$ , let  $\text{Diag } v \in \mathbb{C}^{d \times d}$  be the diagonal matrix whose diagonal entries represent  $v$ . Conversely, for every square matrix  $M$ , let  $\text{diag } M$  be the vector consisting of the diagonal entries of  $M$ . For a matrix (or generally a linear operator)  $M$ , we use  $\|M\|_{\text{Tr}}$  and  $\|M\|_2$  to denote its trace norm and spectral norm respectively, that is,

$$\|M\|_{\text{Tr}} = \text{Tr} \left[ \sqrt{MM^\dagger} \right], \quad \|M\|_2 = \max_{v \neq 0} \|Mv\|_2 / \|v\|_2.$$

For an Hermitian  $M \in \mathbb{C}^{d \times d}$ , we say it is a positive semi-definite operator if for every  $v \in \mathbb{C}^d$ ,  $v^\dagger M v \geq 0$ . A (partial) density operator is a positive semi-definite operator with its trace being 1 (or at most 1, respectively).

#### Viewing a Learning Problem as a Matrix

Let  $M : \mathcal{A} \times \mathcal{X} \rightarrow \{-1, 1\}$  be a matrix. The matrix  $M$  corresponds to the following learning problem. There is an unknown element  $x \in \mathcal{X}$  that was chosen uniformly at random. A learner tries to learn  $x$  from samples  $(a, b)$ , where  $a \in \mathcal{A}$  is chosen uniformly at random and  $b = M(a, x)$ . That is, the learning algorithm is given a stream of samples,  $(a_1, b_1), (a_2, b_2), \dots$ , where each  $a_t$  is uniformly distributed and for every  $t$ ,  $b_t = M(a_t, x)$ . For each  $a \in \mathcal{A}$ , we use  $M_a : \mathcal{X} \rightarrow \{-1, 1\}$  to denote the vector corresponding to the  $a$ -th row of  $M$ .

#### Extractors

A matrix  $M : \mathcal{A} \times \mathcal{X} \rightarrow \{-1, 1\}$  with  $n = \log_2 |\mathcal{X}|$  is a  $(k, \ell)$ - $L_2$  extractor with error  $2^{-r}$ , if for every distribution  $P$  over  $\mathcal{X}$  with  $\|P\|_2 \leq 2^\ell \cdot 2^{-n/2}$ , there are at most  $2^{-k} \cdot |\mathcal{A}|$  rows  $a \in \mathcal{A}$  such that

$$|\langle M_a, P \rangle| \geq 2^{-r}.$$

#### 3.2 Anti-Concentration Bound for Quadratic Form on Unit Vectors

**Lemma 3.1.** *There exists an absolute constant  $c$  such that following holds. Let  $\sigma$  be a Hermitian operator over the Hilbert space  $\mathcal{V} = \mathbb{C}^d$ , and let  $v$  be a uniformly random unit vector in  $\mathcal{V}$ . Then for every  $\varepsilon > 0$ , we have*

$$\Pr \left[ |v^\dagger \sigma v| \leq \frac{\varepsilon \|\sigma\|_2}{d} \right] \leq c\sqrt{\varepsilon} + e^{-d}.$$

*Proof.* Let  $g = (g_1, \dots, g_d) \sim \mathcal{N}(0, 1)^d$  be standard Gaussians. Notice that  $\|g\|_2^2$  follows  $\chi_d^2$ -distribution, and  $|g^\dagger \sigma g| / \|g\|_2^2$  is equidistributed as  $|v^\dagger \sigma v|$ . Therefore by union bound we have

$$\Pr_v \left[ |v^\dagger \sigma v| \leq \frac{\varepsilon \|\sigma\|_2}{d} \right] = \Pr_g \left[ |g^\dagger \sigma g| \leq \varepsilon \|\sigma\|_2 \cdot \frac{\|g\|_2^2}{d} \right] \leq \Pr_g \left[ |g^\dagger \sigma g| \leq 5\varepsilon \|\sigma\|_2 \right] + \Pr_g \left[ \|g\|_2^2 \geq 5d \right].$$

For the first term, notice that  $\text{Var}[g^\dagger \sigma g] = 2\text{Tr}[\sigma^2]$  (see e.g. [RS08, Chapter 5]) which is no smaller than  $2\|\sigma\|_2^2$ . Therefore, by Carbery–Wright inequality [CW01], there exists an absolute constant  $c$  such that

$$\Pr \left[ |g^\dagger \sigma g| \leq 5\varepsilon \|\sigma\|_2 \right] \leq \Pr \left[ |g^\dagger \sigma g| \leq 4\varepsilon \text{Var}[g^\dagger \sigma g]^{1/2} \right] \leq c\sqrt{\varepsilon}.$$

For the second term, the standard Laurent-Massart bound on  $\chi^2$ -distributions [LM00] gives:

$$\Pr \left[ \|g\|_2^2 \geq 5d \right] \leq e^{-d}. \quad \square$$

### 3.3 Multipartite Quantum Systems

The state of  $q$  qubits can be represented in a Hilbert space  $\mathcal{V} = (\mathbb{C}^2)^{\otimes q} = \mathbb{C}^{2^q}$ . In a product of  $m$  Hilbert spaces  $\mathcal{V}_{[m]} = \mathcal{V}_1 \otimes \cdots \otimes \mathcal{V}_m$ , a multipartite partial system  $V_1, \dots, V_m$  is represented by a partial density operator  $\rho_{V_{[m]}}$ . For a subset  $I \subseteq [m]$  of indices, the subsystem on  $\{V_i\}_{i \in I}$  (or  $V_I$  for short) is defined by tracing out  $j \notin I$ , that is,

$$\rho_{V_I} = \text{Tr}_{V_{j \notin I}}[\rho_{V_{[m]}}].$$

Now for any two disjoint subsets  $I, J \subset [m]$ , given some  $|v_J\rangle \in \mathcal{V}_J = \bigotimes_{j \in J} \mathcal{V}_j$ , the conditional system on  $V_I$  is defined as

$$\rho_{V_I|v_J} = (\mathbb{I}_{V_I} \otimes \langle v_J|) \rho_{V_I \cup J} (\mathbb{I}_{V_I} \otimes |v_J\rangle),$$

which is a partial density operator on  $V_I$ . Note that the trace

$$\text{Tr} \left[ \rho_{V_I|v_J} \right] = \langle v_J | \rho_{V_J} | v_J \rangle$$

only depends on the system  $\rho$  and  $|v_J\rangle$ , while being *independent* of the choice of  $I$ .

Another simple fact that will be repeatedly used later on is that for an *orthogonal basis*  $\mathcal{B}$  of  $V_J$ , we have

$$\rho_{V_I} = \text{Tr}_{V_J}[\rho_{V_I \cup J}] = \sum_{|v_J\rangle \in \mathcal{B}} \rho_{V_I|v_J}.$$

### 3.4 Classical-Quantum Systems

In the underlying space  $\mathcal{V}_1 \otimes \cdots \otimes \mathcal{V}_m$  of the multipartite system, we say  $\mathcal{V}_i$  is classical if there is a fixed orthogonal basis  $\mathcal{B}_i$  of  $\mathcal{V}_i$ , such that for every multipartite system  $\rho_{V_{[m]}}$ , every pair of distinct  $|v_i\rangle \neq |v'_i\rangle \in \mathcal{B}_i$  and every two states  $|v\rangle, |v'\rangle \in \bigotimes_{j \neq i} \mathcal{V}_j$ , we have

$$\langle v_i, v | \rho_{V_{[m]}} | v'_i, v' \rangle = 0.$$

Without loss of generality, in the rest of the work we always assume  $\mathcal{B}_i$  is the set of computational basis states. We also identify  $\mathcal{V}_i$  with the discrete set  $\mathcal{B}_i$ , and remove the Dirac brackets when we talk about the classical elements in  $\mathcal{V}_i$ . In this case every multipartite system  $\rho_{V_{[m]}}$  can be written as a direct sum

$$\rho_{V_{[m]}} = \bigoplus_{v_i \in \mathcal{V}_i} \rho_{V_{[m] \setminus \{i\}} | v_i}.$$

The reader may find this direct sum viewpoint easier to handle in some later scenarios.

When  $V_I$  is classical, conditioned on any  $|v_J\rangle \in \mathcal{V}_J$  with  $J$  disjoint from  $I$ , the system  $\rho_{V_I|v_J}$  is represented as a diagonal matrix on  $V_I$ . If  $\text{Tr}[\rho_{V_I|v_J}] > 0$ , it induces a distribution over the computation basis states of  $V_I$ , defined as

$$P_{V_I|v_J}^\rho = \text{diag} \rho_{V_I|v_J} / \text{Tr}[\rho_{V_I|v_J}]. \quad (3)$$

In the rest of this paper, whenever we use this notation  $P_{V_I|v_J}^\rho$ , it is always implicitly assumed that  $\text{Tr}[\rho_{V_I|v_J}] > 0$  and the distribution exists.

In this work we typically consider the following scenario: There is a quantum memory register  $V$  ranging in the Hilbert space  $\mathcal{V}$ , and a classical memory register  $W$  ranging in the set of memory states  $\mathcal{W}$ , along with some classical information  $X \in \mathcal{X}$  (later in the work, it is the concept to be learned) that is correlated with  $V$  and  $W$ . We will make use of the following fact:

**Claim 3.2.** *Let  $\rho_{XVW}$  be a classical-quantum system over classical  $X, W$  and quantum  $V$ . For every  $w \in \mathcal{W}$ ,  $P_{X|w}^\rho$  is a convex combination of  $P_{X|v,w}^\rho$  for some  $\{|v\rangle\} \subseteq \mathcal{V}$ .*

*Proof.* Let  $\mathcal{B}$  be an orthogonal basis of  $\mathcal{V}$ , so that we have (from the end of last section)

$$\rho_{X|w} = \sum_{|v\rangle \in \mathcal{B}} \rho_{X|v,w}.$$

Therefore  $P_{X|w}^\rho$  is a linear combination of  $P_{X|v,w}^\rho$  for  $|v\rangle \in \mathcal{B}$ , with non-negative coefficients. Since they are all distributions, it must be a convex combination.  $\square$

**Characterization of operators over classical-quantum hybrid systems.** Now we identify all possible operators on the classical-quantum hybrid memory space  $\mathcal{V} \otimes \mathcal{W}$ . A priori to the assumption that  $W$  is classical, we think of a quantum channel operating on the system as working on the underlying space  $\mathcal{V} \otimes \mathbb{C}^{|\mathcal{W}|}$ . Now we denote  $\mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$  to be the set of all such quantum channels  $\Phi$  that satisfy the following: for every classical-quantum system  $\rho_{VW}$  in  $\mathcal{V} \otimes \mathcal{W}$ ,  $W$  is still classical in  $\Phi(\rho_{VW})$ . That is, for every two states  $|v\rangle, |v'\rangle \in \mathcal{V}$  and every pair of distinct  $w, w' \in \mathcal{W}$ , we have

$$\langle v, w | \Phi(\rho_{VW}) | v', w' \rangle = 0.$$

Note that not all channels in  $\mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$  are physically realizable. For instance, with one-bit classical memory and no quantum memory, the channel

$$\begin{pmatrix} a & c \\ \bar{c} & b \end{pmatrix} \mapsto \begin{pmatrix} a & ic \\ -i\bar{c} & b \end{pmatrix}$$

is not a classical operator. However, since we are constrained to classical quantum systems, this channel is effectively equivalent to an identity channel on one-bit classical memory. Generally speaking, every channel in  $\mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$  is equivalent to a channel controlled by  $\mathcal{W}$  that maps  $\mathcal{V}$  to  $\mathcal{V} \otimes \mathcal{W}$ . Below, we prove this observation and use it to show the following claim:

**Claim 3.3.** *Let  $\rho_{XVW}$  be a classical-quantum system over classical  $X, W$  and quantum  $V$ . Let  $\Phi \in \mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$ , and we use  $\Phi(\rho)$  to denote the system after applying  $\Phi$  to  $VW$  and identity to  $X$ . Then for every  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$ ,  $P_{X|v,w}^{\Phi(\rho)}$  is a convex combination of  $P_{X|v',w'}^\rho$  for some  $\{|v'\rangle\} \subseteq \mathcal{V}$  and  $\{w'\} \subseteq \mathcal{W}$ .*

One difference between Claim 3.2 and Claim 3.3 is that in Claim 3.3 it is not always possible to write  $P_{X|v,w}^{\Phi(\rho)}$  as a convex combination of  $P_{X|v',w'}^\rho$  for  $|v'\rangle$  from an orthogonal basis of  $\mathcal{V}$ . But it is always possible in Claim 3.2. Although the difference does not matter in this work, we mention it here for clarity.

*Proof.* Since  $\Phi \in \mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$ , the following channel is functionally equivalent to  $\Phi$  for classical-quantum systems:

$$\Phi' : \rho \rightarrow \sum_{w \in \mathcal{W}} \Phi(\rho_{V|w} \otimes |w\rangle\langle w|).$$

The physical meaning of  $\Phi'$  is to measure  $W$  under the computational basis (which should not change the functionality we care about) and apply  $\Phi$ .

By defining the channel  $\Phi_w(\cdot) := \Phi(\cdot \otimes |w\rangle\langle w|)$ , the above can be alternatively written as:

$$\Phi' : \rho \rightarrow \sum_{w \in \mathcal{W}} \Phi_w(\rho_{V|w}).$$

Now consider the Kraus representation of each  $\Phi_w$ , that is, a finite set of linear operators  $E_{w,k} : \mathcal{V} \rightarrow \mathcal{V} \otimes \mathcal{W}$  such that

$$\Phi_w(\rho_{V|w}) = \sum_k E_{w,k} \rho_{V|w} E_{w,k}^\dagger, \quad \sum_k E_{w,k}^\dagger E_{w,k} = \mathbb{I}_V.$$

We can write

$$\begin{aligned} \Phi(\rho)_{X|v,w} &= \Phi'(\rho)_{X|v,w} = (\mathbb{I}_X \otimes \langle v, w |) \Phi'(\rho) (\mathbb{I}_X \otimes |v, w \rangle) \\ &= \sum_{w' \in \mathcal{W}} \sum_k (\mathbb{I}_X \otimes \langle v, w | E_{w',k}) \rho_{XV|w'} (\mathbb{I}_X \otimes E_{w',k}^\dagger |v, w \rangle) \\ &= \sum_{w' \in \mathcal{W}} \sum_k \|E_{w',k}^\dagger |v, w \rangle\|_2 \cdot \rho_{X|v',w'} \end{aligned}$$

where in each term of the summation,  $|v'\rangle \sim E_{w',k}^\dagger |v, w\rangle$ . Similar to the arguments in Claim 3.2,  $P_{X|v,w}^{\Phi(\rho)}$  is a convex combination of  $P_{X|v',w'}^\rho$ .  $\square$

### 3.5 Branching Program with Hybrid Memory

For a learning problem that corresponds to the matrix  $M$ , a branching program of hybrid memory with  $m$ -bit classical memory,  $q$ -qubit quantum memory and length  $T$  is specified as follows.

At each stage  $0 \leq t \leq T$ , the memory state of the branching program is described as a classical-quantum system  $\rho_{VW}^{(t)}$  over quantum memory space  $\mathcal{V} = (\mathbb{C}^2)^{\otimes q}$  and classical memory space  $\mathcal{W} = \{0, 1\}^m$ . The memory state evolves based on the samples that the branching program receives, and therefore depends on the unknown element  $x \in_R \mathcal{X}$ . We can then interpret the overall systems over  $XVW$ , in which  $X$  consists of an unknown concept  $x$ , resulting in a classical-quantum system  $\rho_{XVW}^{(t)}$ . It always holds that the distribution of  $x$  is uniform, i.e.,

$$\rho_X^{(t)} = \text{Tr}_{VW}[\rho_{XVW}^{(t)}] = \frac{1}{2^n} \mathbb{I}_X.$$

Initially the memory  $VW$  is independent of  $X$  and can be arbitrarily initialized. We assume that

$$\rho_{XVW}^{(0)} = \frac{1}{2^n} \mathbb{I}_X \otimes \frac{1}{2^q} \mathbb{I}_V \otimes \frac{1}{2^m} \mathbb{I}_W.$$

At each stage  $0 \leq t < T$ , the branching program receives a sample  $(a, b)$ , where  $a \in_R \mathcal{A}$  and  $b = M(a, x)$ , and applies an operation  $\Phi_{t,a,b} \in \mathcal{T}_{\mathcal{V} \otimes \mathcal{W}}$  over its memory state. Thus the evolution of the entire system can be written as

$$\rho_{XVW}^{(t+1)} = \mathbf{E}_{a \in_R \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t,a,M(a,x)}(\rho_{VW|x}^{(t)}) \right].$$

Finally, at stage  $t = T$ , a measurement over the computational bases is applied on  $\rho_{VW}^{(T)}$ , and the branching program outputs an element  $\tilde{x} \in X$  as a function of the measurement result  $(v, w) \in \{0, 1\}^{q+m}$ . The success probability of the program is the probability that  $\tilde{x} = x$  which can be formulated as

$$\sum_{\substack{x \in \mathcal{X}, v \in \{0,1\}^q, w \in \mathcal{W} \\ \tilde{x}(v,w)=x}} \langle x, v, w | \rho_{XVW}^{(T)} | x, v, w \rangle.$$

## 4 Main Result

**Theorem 2.** *Let  $\mathcal{X}, \mathcal{A}$  be two finite sets with  $n = \log_2 |\mathcal{X}|$ . Let  $M : \mathcal{A} \times \mathcal{X} \rightarrow \{-1, 1\}$  be a matrix which is a  $(k', \ell')$ - $L_2$  extractor with error  $2^{-r'}$  for sufficiently large  $k', \ell'$  and  $r'$ , where  $\ell' \leq n$ . Let*

$$r = \min \left\{ \frac{1}{4}r', \frac{1}{26}\ell' + \frac{1}{6}, \frac{1}{2}(k' - 1) \right\}.$$

*Let  $\rho$  be a branching program for the learning problem corresponding to  $M$ , described by classical-quantum systems  $\rho_{XVW}^{(t)}$ , with  $q$ -qubit quantum memory  $V$ ,  $m$ -bit classical memory  $W$  and length  $T$ . If  $m \leq \frac{1}{44}(k' - 1)\ell'$ ,  $q \leq r - 7$  and  $T \leq 2^{r-2}$ , the success probability of  $\rho$  is at most  $O(2^{q-r})$ .*

From now on we let  $k = k' - 1$  and  $\ell = \frac{1}{5}(\ell' - 13r - 2)$ . Then we have the following inequalities to be used later:

$$q + r + 1 - r' \leq -2r. \quad (4)$$

$$2\ell + 9r - n \leq -r. \quad (5)$$

$$(k - r)\ell \geq 2m + 4r + 1. \quad (6)$$

We leave a detailed calculation for the above inequalities in Appendix A.

### 4.1 Truncated Classical-Quantum Systems

Here we describe how to truncate a partial classical-quantum system  $\rho_{XVW}$  according to some property  $G(v, w)$  of desire on  $\rho_{X|v, w}$ . The goal is to remove the parts of  $\rho_{XVW}$  where  $G$  is not satisfied. We execute the following procedure:

1. Maintain a partial system  $\rho'_{XVW}$  initialized as  $\rho_{XVW}$ , and subspaces  $\mathcal{V}_w \subseteq \mathcal{V}$  initialized as  $\mathcal{V}$  for each  $w \in \mathcal{W}$ .
2. Pick  $w \in \mathcal{W}$  and  $|v\rangle \in \mathcal{V}_w$  such that  $\text{Tr}[\rho'_{X|v, w}] > 0$  and  $G(v, w)$  is false.
3. Change the partial system  $\rho'_{XVW}$  into the following system by projection:

$$(\mathbb{I}_X \otimes (\mathbb{I}_{VW} - |v, w\rangle\langle v, w|))\rho'_{XVW}(\mathbb{I}_X \otimes (\mathbb{I}_{VW} - |v, w\rangle\langle v, w|)),$$

and change  $\mathcal{V}_w$  to its subspace orthogonal to  $|v\rangle$ , that is

$$\{|v'\rangle \in \mathcal{V}_w \mid \langle v|v'\rangle = 0\}.$$

4. Repeat from step 2 until there is no such  $w$  and  $|v\rangle$ . Denote the final system as  $\rho_{XVW}^{|G}$ .

In step 2 we pick  $w$  and  $|v\rangle$  arbitrarily as long as it satisfies the requirements, however we could always think of it as iterating over  $w \in \mathcal{W}$  and processing each  $\rho_{XV|w}$  separately. The choices of  $|v\rangle$  for each  $w$  do affect the final system  $\rho_{XVW}^{|G}$ ; Yet as we will see later, these choices are irrelevant to our proof.

Below, we give two useful lemmas on truncated systems.

**Lemma 4.1.** *For every  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$  such that  $\text{Tr}[\rho_{X|v, w}^{|G}] > 0$ , there exists  $|v'\rangle$  in the remaining subspace  $\mathcal{V}_w$  such that*

$$P_{X|v, w}^{\rho^{|G}} = P_{X|v', w}^{\rho} = P_{X|v', w}^{\rho^{|G}}.$$

*Proof.* It suffices to prove the lemma with one round of the truncation procedure executed. Suppose the  $|v_1, w_1\rangle$  is picked in step 2, resulting in the partial system

$$\rho'_{XVW} = (\mathbb{I}_X \otimes (\mathbb{I}_{VW} - |v_1, w_1\rangle\langle v_1, w_1|))\rho_{XVW}(\mathbb{I}_X \otimes (\mathbb{I}_{VW} - |v_1, w_1\rangle\langle v_1, w_1|)).$$

We can write

$$\begin{aligned}\rho'_{X|v,w} &= (\mathbb{I}_X \otimes \langle v, w|)\rho'_{XVW}(\mathbb{I}_X \otimes |v, w\rangle) \\ &= (\mathbb{I}_X \otimes (\langle v, w| - \langle v, w|v_1, w_1\rangle\langle v_1, w_1|))\rho_{XVW}(\mathbb{I}_X \otimes (|v, w\rangle - |v_1, w_1\rangle\langle v_1, w_1|v, w\rangle)).\end{aligned}$$

- If  $w \neq w_1$ , then

$$\rho'_{X|v,w} = (\mathbb{I}_X \otimes \langle v, w|)\rho_{XVW}(\mathbb{I}_X \otimes |v, w\rangle) = \rho_{X|v,w}.$$

And the lemma holds directly by choosing  $|v'\rangle = |v\rangle$ .

- If  $w = w_1$ , then with  $\langle v_1, w_1|v, w\rangle = \langle v_1|v\rangle = \lambda$ , we have

$$\rho'_{X|v,w} = (\mathbb{I}_X \otimes (\langle v| - \bar{\lambda}\langle v_1|)\langle w|)\rho_{XVW}(\mathbb{I}_X \otimes (|v\rangle - \lambda|v_1\rangle)|w\rangle).$$

By the fact that  $\text{Tr}[\rho_{X|v,w}^{|G|}] > 0$ , we must have  $|v\rangle \neq |v_1\rangle$ . Therefore if we let  $|v'\rangle \sim |v\rangle - \lambda|v_1\rangle$ , which is the normalized projection of  $|v\rangle$  onto the orthogonal subspace of  $|v_1\rangle$ , the above equality implies that  $P_{X|v,w}^{\rho'} = P_{X|v',w}^{\rho}$ . Meanwhile, since  $\langle v_1|v'\rangle = 0$  we have  $\rho'_{X|v',w} = \rho_{X|v',w}$ , which completes the proof.  $\square$

A direct corollary of the above lemma is that if  $G(v, w)$  only depends on the distribution  $P_{X|v,w}^{\rho}$ , then  $G(v, w)$  holds for every  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$  in the truncated system  $\rho_{XVW}^{|G|}$ , even when  $|v\rangle$  is not in the remaining subspace  $\mathcal{V}_w$ .

Our second lemma is based on the following fact that bounds the trace distance of a partial system and its projection, whose proof can be found in the Appendix B.

**Proposition 4.2.** *For every partial system  $\rho$  and projection operator  $\Pi$  on  $\rho$ , we have*

$$\|\rho - \Pi\rho\Pi\|_{\text{Tr}}^2 \leq 4\text{Tr}[\rho]^2 - 4\text{Tr}[\Pi\rho]^2.$$

**Lemma 4.3.** *For each  $w \in W$ , let  $|v_1\rangle, \dots, |v_d\rangle$  be the states picked in step 2 within  $\mathcal{V}_w$ . Then*

$$\|\rho_{XV|w} - \rho_{XV|w}^{|G|}\|_{\text{Tr}} \leq 3 \sum_{i=1}^d \sqrt{\text{Tr}[\rho_{X|v_i,w}] \text{Tr}[\rho_{XV|w}]}.$$

*Proof.* In Proposition 4.2, take  $\rho$  to be  $\rho_{XV|w}$ , and  $\Pi$  to be

$$\mathbb{I}_X \otimes \prod_{i=1}^d (\mathbb{I}_V - |v_i\rangle\langle v_i|) = \mathbb{I}_X \otimes \left( \mathbb{I}_V - \sum_{i=1}^d |v_i\rangle\langle v_i| \right).$$

Then  $\Pi\rho\Pi = \rho_{XV|w}^{|G|}$  and  $\text{Tr}[\Pi\rho] = \text{Tr}[\rho_{XV|w}] - \sum_{i=1}^d \text{Tr}[\rho_{X|v_i,w}]$ . Therefore we have

$$\begin{aligned}\|\rho_{XV|w} - \rho_{XV|w}^{|G|}\|_{\text{Tr}} &\leq \sqrt{4\text{Tr}[\rho]^2 - 4\text{Tr}[\Pi\rho]^2} \\ &\leq \sqrt{8(\text{Tr}[\rho] - \text{Tr}[\Pi\rho])\text{Tr}[\rho]} \\ &= \sqrt{8 \sum_{i=1}^d \text{Tr}[\rho_{X|v_i,w}] \text{Tr}[\rho_{XV|w}]} \\ &\leq 3 \sum_{i=1}^d \sqrt{\text{Tr}[\rho_{X|v_i,w}] \text{Tr}[\rho_{XV|w}]}.\end{aligned}\quad \square$$



Since  $\text{Tr}[\rho_{XV|w}] \leq 1$  always holds, by summing over all  $w \in \mathcal{W}$  we get the following corollary:

**Corollary 4.4.** *Let  $|v_1, w_1\rangle, \dots, |v_d, w_d\rangle$  be all of the memory states picked in step 2. Then*

$$\|\rho_{XVW} - \rho_{XVW}^G\|_{\text{Tr}} \leq 3 \sum_{i=1}^d \sqrt{\text{Tr}[\rho_{X|v_i, w_i}]}$$

## 4.2 Truncated Branching Program

The properties that we desire for the partial system  $\rho_{XVW}$  consist of three parts:

- Small  $L_2$  norm: Let  $G_2(v, w)$  be the property that

$$\|P_{X|v, w}^\rho\|_2 \leq 2^\ell \cdot 2^{-n/2}.$$

- Small  $L_\infty$  norm: Let  $G_\infty(v, w)$  be the property that

$$\|P_{X|v, w}^\rho\|_\infty \leq 2^{2\ell+9r} \cdot 2^{-n}.$$

- Even division: For every  $a \in \mathcal{A}$ , let  $G_a(v, w)$  be the property that

$$|\langle M_a, P_{X|v, w}^\rho \rangle| \leq 2^{-r}.$$

Now we define the truncated branching program, by specifying the truncated partial classical-quantum system  $\tau_{XVW}^{(t)}$  for each stage  $t$ . Initially let  $\tau_{XVW}^{(0)} = \rho_{XVW}$ . For each stage  $0 \leq t \leq T$ , the truncation consists of three ingredients (below we ignore the superscripts on  $P$  for convenience):

1. Remove parts where  $\|P_{X|v, w}\|_2$  is large. That is, let  $\tau_{XVW}^{(t, \star)} = \tau_{XVW}^{(t)|G_2}$ .
2. Remove parts where  $\|P_{X|v, w}\|_\infty$  is large. This is done by two steps.

- First, let  $g \in \{0, 1\}^{\mathcal{X} \otimes \mathcal{W}}$  be an indicator vector such that  $g(x, w) = 1$  if and only if

$$\text{Tr}[\tau_{X|w}^{(t, \star)}] > 0 \text{ and } P_{X|w}^{\tau_{X|w}^{(t, \star)}}(x) \leq 2^{2\ell+5r} \cdot 2^{-n}.$$

Let  $\tau_{XVW}^{(t, \circ)} = (gg^\dagger \otimes \mathbb{I}_V) \tau_{XVW}^{(t, \star)} (gg^\dagger \otimes \mathbb{I}_V)$ , where  $gg^\dagger$  is the projection operator acting on  $\mathcal{X} \otimes \mathcal{W}$ .

- To make sure that the distributions did not change a lot after the projection  $gg^\dagger$ , for each  $0 \leq t < T$ , let  $G_t(v, w)$  be the property that

$$\text{Tr}[\tau_{X|v, w}^{(t, \circ)}] \geq (1 - 2^{-r}) \text{Tr}[\tau_{X|v, w}^{(t, \star)}].$$

Let  $\tau_{XVW}^{(t, \infty)} = \tau_{XVW}^{(t, \circ)|G_\infty \wedge G_t}$ .

3. For each  $a \in \mathcal{A}$ , remove (only for this  $a$ ) parts where  $P_{X|v, w}$  is not evenly divided by  $a$ . That is, for each  $a \in \mathcal{A}$ , let  $\tau_{XVW}^{(t, a)} = \tau_{XVW}^{(t, \infty)|G_a}$ .

Then, if  $t < T$ , for each  $a \in_R \mathcal{A}$  we evolve the system by applying the sample operations  $\Phi_{t, a, b}$  as the original branching program on  $\tau_{XVW}^{(t, a)}$ , so that we have

$$\tau_{XVW}^{(t+1)} = \mathbf{E}_{a \in_R \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t, a, M(a, x)}(\tau_{VW|x}^{(t, a)}) \right].$$

### 4.3 Bounding the Truncation Difference

In order to show that the success probability of the original branching program  $\rho^{(t)}$  is low, the plan is to prove an upper bound on the success probability of the truncated branching program  $\tau^{(t)}$ , and bound the difference between the two probabilities.

Here we bound the difference by the trace distance between the two systems  $\rho_{XVW}^{(t)}$  and  $\tau_{XVW}^{(t)}$ . We will show that the contribution to the trace distance from each one of the truncation ingredients is small, and in addition the evolution preserves the trace distance.

#### 4.3.1 Truncation by $G_2$

**Lemma 4.5.** *For every  $0 \leq t \leq T$ ,  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$  such that  $G_2(v, w)$  is violated (that is,  $\|P_X^{\tau^{(t)}}\|_2 > 2^\ell \cdot 2^{-n/2}$ ), we must have  $\text{Tr}[\tau_{X|v,w}^{(t)}] < 2^{-2m} \cdot 2^{-4r}$ .*

The lemma says, for any direction  $|v, w\rangle$  picked by the truncation procedure, the weight will be small and the truncation will not change the state significantly.

*Proof.* This is our main technical lemma and we defer the proof to Section 5.  $\square$

Since there are at most  $2^{q+m}$  such directions picked in the truncation procedure, we conclude the following corollary.

**Corollary 4.6.** *For every  $0 \leq t \leq T$ , we have  $\|\tau_{XVW}^{(t,*)} - \tau_{XVW}^{(t)}\|_{\text{Tr}} \leq 3 \cdot 2^{q-2r}$ .*

*Proof.* Recall that  $\tau_{XVW}^{(t,*)} = \tau_{XVW}^{(t)|G_2}$ . Since  $\dim(\mathcal{V} \otimes \mathcal{W}) = 2^{q+m}$ , the truncation lasts for at most  $2^{q+m}$  rounds. Since in every round the picked  $|v, w\rangle$  has  $\text{Tr}[\tau_{X|v,w}^{(t)}] < 2^{-2m} \cdot 2^{-4r}$ , by Corollary 4.4 we have

$$\|\tau_{XVW}^{(t,*)} - \tau_{XVW}^{(t)}\|_{\text{Tr}} \leq 3 \cdot 2^{q+m} \cdot \sqrt{2^{-2m} \cdot 2^{-4r}} = 3 \cdot 2^{q-2r}. \quad \square$$

#### 4.3.2 Truncation by $G_\infty$

**Lemma 4.7.** *For every  $0 \leq t \leq T$  and  $w \in \mathcal{W}$  we have*

$$\sum_{\substack{x \in \mathcal{X} \\ g(x,w)=0}} P_X^{\tau^{(t,*)}}(x) \leq 2^{-5r}.$$

*Proof.* By Claim 3.2,  $P_X^{\tau^{(t,*)}}$  is a convex combination of  $P_X^{\tau^{(t,*)}}_{|v,w}$ . From Lemma 4.1 we know that  $G_2(P_X^{\tau^{(t,*)}}_{|v,w})$  holds for every  $|v\rangle$  and  $w$ , and thus by convexity of  $\ell_2$ -norms we know that  $G_2(P_X^{\tau^{(t,*)}})$  also holds. That means

$$\mathbf{E}_{x \sim P_X^{\tau^{(t,*)}}} \left[ P_X^{\tau^{(t,*)}}(x) \right] = \|P_X^{\tau^{(t,*)}}\|_2^2 \leq 2^{2\ell} \cdot 2^{-n}.$$

Therefore, by Markov's inequality we have

$$\sum_{\substack{x \in \mathcal{X} \\ g(x,w)=0}} P_X^{\tau^{(t,*)}}(x) = \Pr_{x \sim P_X^{\tau^{(t,*)}}} \left[ P_X^{\tau^{(t,*)}}(x) > 2^{2\ell+5r} \cdot 2^{-n} \right] \leq 2^{-5r}. \quad \square$$

**Corollary 4.8.** *For every  $0 \leq t \leq T$  and every  $w \in \mathcal{W}$ , we have  $\tau_{XV|w}^{(t,\circ)} \leq \tau_{XV|w}^{(t,*)}$ , and*

$$\text{Tr}[\tau_{XV|w}^{(t,\circ)}] \geq (1 - 2^{-5r}) \cdot \text{Tr}[\tau_{XV|w}^{(t,*)}].$$

Moreover, we have  $\|\tau_{XVW}^{(t,\circ)} - \tau_{XVW}^{(t,*)}\|_{\text{Tr}} \leq 2^{-5r}$ .

*Proof.* Since  $X$  and  $W$  are both classical and  $\tau_{XVW}^{(t,\circ)} = (gg^\dagger \otimes \mathbb{I}_V) \tau_{XVW}^{(t,\star)} (gg^\dagger \otimes \mathbb{I}_V)$ , we have

$$\tau_{XV|w}^{(t,\star)} - \tau_{XV|w}^{(t,\circ)} = \sum_{\substack{x \in \mathcal{X} \\ g(x,w)=0}} |x\rangle\langle x| \otimes \tau_{V|x,w}^{(t,\star)},$$

which is positive semi-definite. Recalling that (Equation (3))

$$\mathrm{Tr}[\tau_{V|x,w}^{(t,\star)}] = \langle x, w | \tau_{XW}^{(t,\star)} |x, w\rangle = \mathrm{diag} \tau_{X|w}^{(t,\star)}(x) = P_{X|w}^{\tau^{(t,\star)}}(x) \mathrm{Tr}[\tau_{X|w}^{(t,\star)}],$$

we have

$$\mathrm{Tr}[\tau_{XV|w}^{(t,\star)}] - \mathrm{Tr}[\tau_{XV|w}^{(t,\circ)}] = \sum_{\substack{x \in \mathcal{X} \\ g(x,w)=0}} \mathrm{Tr}[\tau_{V|x,w}^{(t,\star)}] = \sum_{\substack{x \in \mathcal{X} \\ g(x,w)=0}} P_{X|w}^{\tau^{(t,\star)}}(x) \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\star)}] \leq 2^{-5r} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\star)}].$$

And therefore, as  $\tau_{XVW}^{(t,\circ)} - \tau_{XVW}^{(t,\star)}$  is positive semi-definite, we have

$$\|\tau_{XVW}^{(t,\circ)} - \tau_{XVW}^{(t,\star)}\|_{\mathrm{Tr}} = \sum_{w \in \mathcal{W}} \mathrm{Tr}[\tau_{XV|w}^{(t,\star)}] - \mathrm{Tr}[\tau_{XV|w}^{(t,\circ)}] \leq 2^{-5r} \sum_{w \in \mathcal{W}} \mathrm{Tr}[\tau_{X|w}^{(t,\star)}] \leq 2^{-5r}. \quad \square$$

**Lemma 4.9.** *For every  $0 \leq t \leq T$ ,  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$  such that  $G_\infty(v, w)$  is violated (that is,  $\|P_{X|v,w}^{\tau^{(t,\circ)}}\|_\infty > 2^{2\ell+9r} \cdot 2^{-n}$ ) or  $G_t(v, w)$  is violated (that is,  $\mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] < (1 - 2^{-r})\mathrm{Tr}[\tau_{X|v,w}^{(t,\star)}]$ ), we must have  $\mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] < 2 \cdot 2^{-4r} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}]$ .*

*Proof.* If  $G_\infty(v, w)$  is violated, let  $x \in \mathcal{X}$  be the one such that  $P_{X|v,w}^{\tau^{(t,\circ)}}(x) > 2^{2\ell+9r} \cdot 2^{-n}$ . If  $g(x, w) = 0$  then  $P_{X|w}^{\tau^{(t,\circ)}}(x) = 0$ , while if  $g(x, w) = 1$  then by Corollary 4.8,

$$P_{X|w}^{\tau^{(t,\circ)}}(x) \leq \frac{\mathrm{Tr}[\tau_{X|w}^{(t,\star)}]}{\mathrm{Tr}[\tau_{X|w}^{(t,\circ)}]} \cdot 2^{2\ell+5r} \cdot 2^{-n} \leq (1 - 2^{-5r})^{-1} \cdot 2^{2\ell+5r} \cdot 2^{-n}.$$

Hence we always have

$$\mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] \leq \frac{P_{X|w}^{\tau^{(t,\circ)}}(x)}{P_{X|v,w}^{\tau^{(t,\circ)}}(x)} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}] \leq 2 \cdot 2^{-4r} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}],$$

where the first inequality comes from the fact that  $\tau_{X|w}^{(t,\circ)} \geq \tau_{X|v,w}^{(t,\circ)}$  and Equation (3).

If  $G_t(v, w)$  is violated, since we know from Corollary 4.8 that

$$\begin{aligned} \left| \mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] - \mathrm{Tr}[\tau_{X|v,w}^{(t,\star)}] \right| &\leq \|\tau_{XV|w}^{(t,\circ)} - \tau_{XV|w}^{(t,\star)}\|_{\mathrm{Tr}} \leq 2^{-5r} \cdot \mathrm{Tr}[\tau_{XV|w}^{(t,\star)}] \\ &\leq 2^{-5r} \cdot (1 - 2^{-5r})^{-1} \cdot \mathrm{Tr}[\tau_{XV|w}^{(t,\circ)}], \end{aligned}$$

therefore from  $\mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] < (1 - 2^{-r})\mathrm{Tr}[\tau_{X|v,w}^{(t,\star)}]$  we deduce that

$$\begin{aligned} \mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] &< (2^r - 1) \cdot \left( \mathrm{Tr}[\tau_{X|v,w}^{(t,\star)}] - \mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] \right) \\ &\leq (2^r - 1) \cdot 2^{-5r} \cdot (1 - 2^{-5r})^{-1} \cdot \mathrm{Tr}[\tau_{XV|w}^{(t,\circ)}] \\ &< 2 \cdot 2^{-4r} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}]. \quad \square \end{aligned}$$

**Corollary 4.10.** *For every  $0 \leq t \leq T$ , we have  $\|\tau_{XVW}^{(t,\infty)} - \tau_{XVW}^{(t,\circ)}\|_{\mathrm{Tr}} \leq 5 \cdot 2^{q-2r}$ .*

*Proof.* Recall that  $\tau_{XVW}^{(t,\infty)} = \tau_{XVW}^{(t,\circ)G_\infty \wedge G_t}$ . For each  $w \in \mathcal{W}$ , the truncation picks at most  $\dim \mathcal{V} = 2^q$  states  $|v, w\rangle$ , each with  $\mathrm{Tr}[\tau_{X|v,w}^{(t,\circ)}] < 2 \cdot 2^{-4r} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}]$ . Therefore by applying Lemma 4.3 for each  $w \in \mathcal{W}$ , we have

$$\|\tau_{XVW}^{(t,\infty)} - \tau_{XVW}^{(t,\circ)}\|_{\mathrm{Tr}} \leq 3 \cdot \sum_{w \in \mathcal{W}} 2^q \cdot \sqrt{2 \cdot 2^{-4r}} \cdot \mathrm{Tr}[\tau_{X|w}^{(t,\circ)}] \leq 5 \cdot 2^{q-2r}. \quad \square$$

### 4.3.3 Truncation by $G_a$

Notice that in the truncation step from  $\tau^{(t,\star)}$  to  $\tau^{(t,\circ)}$ , the distribution  $P_{X|v,w}^{\tau^{(t,\star)}}$  might change and not satisfy  $G_2$  anymore. However, with the truncation by  $G_t$ , any such distribution that changes too much is eliminated, and we have the following guarantee.

**Lemma 4.11.** *For every  $0 \leq t \leq T$ ,  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$ , we have*

$$\|P_{X|v,w}^{\tau^{(t,\infty)}}\|_2 \leq (1 - 2^{-r})^{-1} \cdot 2^\ell \cdot 2^{-n/2}.$$

*Proof.* By Lemma 4.1, there exists  $|v'\rangle \in \mathcal{V}$  such that  $P_{X|v,w}^{\tau^{(t,\infty)}} = P_{X|v',w}^{\tau^{(t,\infty)}} = P_{X|v',w}^{\tau^{(t,\circ)}}$ . The truncation by  $G_t$  ensures that  $\text{Tr}[\tau_{X|v',w}^{(t,\circ)}] \geq (1 - 2^{-r})\text{Tr}[\tau_{X|v',w}^{(t,\star)}]$ , and therefore

$$\|P_{X|v,w}^{\tau^{(t,\infty)}}\|_2 = \|P_{X|v',w}^{\tau^{(t,\circ)}}\|_2 = \frac{\|\text{diag } \tau_{X|v',w}^{(t,\circ)}\|_2}{\text{Tr}[\tau_{X|v',w}^{(t,\circ)}]} \leq \frac{\|\text{diag } \tau_{X|v',w}^{(t,\star)}\|_2}{(1 - 2^{-r})\text{Tr}[\tau_{X|v',w}^{(t,\star)}]} \leq (1 - 2^{-r})^{-1} \cdot 2^\ell \cdot 2^{-n/2}. \quad \square$$

**Lemma 4.12.** *For every partial classical-quantum system  $\tau_{XV}$  over  $\mathcal{X} \otimes \mathcal{V}$  such that  $\|P_{X|v}^\tau\|_2 \leq 2^{\ell'} \cdot 2^{-n/2}$  holds for every  $|v\rangle \in \mathcal{V}$ , we have*

$$\Pr_{a \in_R \mathcal{A}} \left[ \exists |v\rangle \in \mathcal{V}, |\langle M_a, P_{X|v}^\tau \rangle| \geq 2^{-r} \right] \leq 2^{-2r}.$$

*Proof.* Notice that we can think of  $\tau_V = \text{Tr}_X[\tau_{XV}]$  to be  $\mathbb{I}_V$ . This is because we can first assume that  $\tau_V$  is full rank (otherwise change  $\mathcal{V}$  to its subspace and the conclusion in this lemma still holds), and if we have diagonalization  $Q^\dagger \tau_V Q = \mathbb{I}_V$  for some non-singular  $Q$ , then consider the new system

$$\tau'_{XV} = (\mathbb{I}_X \otimes Q^\dagger) \tau_{XV} (\mathbb{I}_X \otimes Q),$$

and the set of distributions  $\{P_{X|v}^\tau\}$  and  $\{P_{X|v}^{\tau'}\}$  over  $|v\rangle \in \mathcal{V}$  are the same, since  $P_{X|v}^{\tau'} = P_{X|v}^\tau$  for  $|v'\rangle \sim Q|v\rangle$ . With  $\tau_V = \mathbb{I}_V$  we have  $\text{Tr}[\tau_{X|v}] = 1$  for every  $|v\rangle \in \mathcal{V}$ , and thus  $P_{X|v}^\tau = \text{diag } \tau_{X|v}$ .

Let  $\mathcal{A}' \subseteq \mathcal{A}$  be the set of  $a \in \mathcal{A}$  such that there exists  $|v\rangle \in \mathcal{V}$  with  $|\langle M_a, P_{X|v}^\tau \rangle| \geq 2^{-r}$ . For each  $a \in \mathcal{A}'$ , let

$$\sigma_a = \text{Tr}_X[(\text{Diag } M_a \otimes \mathbb{I}_V) \tau_{XV}]$$

which is a Hermitian operator on  $\mathcal{V}$ . There exists  $|v\rangle \in \mathcal{V}$  such that

$$|\langle v | \sigma_a | v \rangle| = |\langle M_a, \text{diag } \tau_{X|v} \rangle| = |\langle M_a, P_{X|v}^\tau \rangle| \geq 2^{-r},$$

which means that  $\|\sigma_a\|_2 \geq 2^{-r}$ . Now let  $|u\rangle$  be a uniformly random unit vector in  $\mathcal{V}$ , and by Lemma 3.1 we know that for some absolute constant  $c$ ,

$$\Pr_{|u\rangle} \left[ |\langle u | \sigma_a | u \rangle| \geq 2^{-r'} \right] \geq 1 - 2^{(q+r-r')/2} c - e^{-2q} \geq 1 - 2^{-r} c - e^{-1} \geq 1/2.$$

The second last inequality comes from Eq. (4), while the last inequality is because of the assumption that  $r$  is sufficiently large.

Since the above holds for every  $a \in \mathcal{A}'$ , it implies that  $\Pr_{a \in \mathcal{A}', |u\rangle} [|\langle u | \sigma_a | u \rangle| \geq 2^{-r'}]$  is at least  $1/2$ . It means that there exists some  $|u\rangle \in \mathcal{V}$  such that  $|\langle u | \sigma_a | u \rangle| \geq 2^{-r'}$  for at least half of  $a \in \mathcal{A}'$ . On the other hand, since  $M$  is a  $(k', \ell')$ -extractor with error  $2^{-r'}$ , and  $\|P_{X|u}^\tau\|_2 \leq 2^{\ell'} \cdot 2^{-n/2}$ , there are at most  $2^{-k'}$  fraction of  $a \in \mathcal{A}$  such that  $|\langle u | \sigma_a | u \rangle| = |\langle M_a, P_{X|u}^\tau \rangle| \geq 2^{-r'}$ . That means

$$\Pr_{a \in_R \mathcal{A}} [a \in \mathcal{A}'] \leq 2 \cdot 2^{-k'} \leq 2^{-2r}.$$

Here  $k' - 1 \geq 2r$ , by the definition of  $r$ . □

**Corollary 4.13.** For every  $0 \leq t \leq T$ , we have  $\mathbf{E}_{a \in_R \mathcal{A}} \|\tau_{XVW}^{(t,a)} - \tau_{XVW}^{(t,\infty)}\|_{\text{Tr}} \leq 2^{-2r}$ .

*Proof.* For each  $w \in \mathcal{W}$ , the partial system  $\tau_{XV|w}^{(t,\infty)}$  satisfies the condition of Lemma 4.12 since for every  $|v\rangle \in \mathcal{V}$ ,

$$\|P_X^{\tau_{XV|w}^{(t,\infty)}}\|_2 \leq (1 - 2^{-r})^{-1} \cdot 2^\ell \cdot 2^{-n/2} \leq 2^{\ell'} \cdot 2^{-n/2}.$$

Notice that for each  $a \in \mathcal{A}$  such that there does not exist  $|v\rangle \in \mathcal{V}$  with  $|\langle M_a, P_X^{\tau_{XV|w}^{(t,\infty)}} \rangle| \geq 2^{-r}$  (that is, when  $G_a(v, w)$  holds for every  $|v\rangle \in \mathcal{V}$ ), the sub system  $\tau_{XV|w}^{(t,\infty)}$  is not touched in the truncation by  $G_a$  and we have  $\tau_{XV|w}^{(t,a)} = \tau_{XV|w}^{(t,\infty)}$ . Therefore

$$\begin{aligned} \mathbf{E}_{a \in_R \mathcal{A}} \|\tau_{XVW}^{(t,a)} - \tau_{XVW}^{(t,\infty)}\|_{\text{Tr}} &= \sum_{w \in \mathcal{W}} \mathbf{E}_{a \in_R \mathcal{A}} \|\tau_{XV|w}^{(t,a)} - \tau_{XV|w}^{(t,\infty)}\|_{\text{Tr}} \\ &\leq \sum_{w \in \mathcal{W}} \Pr_{a \in_R \mathcal{A}} \left[ \exists |v\rangle \in \mathcal{V}, |\langle M_a, P_X^{\tau_{XV|w}^{(t,\infty)}} \rangle| \geq 2^{-r} \right] \cdot \text{Tr}[\tau_{XV|w}^{(t,\infty)}] \\ &\leq 2^{-2r} \cdot \sum_{w \in \mathcal{W}} \text{Tr}[\tau_{XV|w}^{(t,\infty)}] \leq 2^{-2r}. \quad \square \end{aligned}$$

#### 4.3.4 Evolution preserves trace distance

**Lemma 4.14.** For every  $0 \leq t < T$ , we have  $\|\tau_{XVW}^{(t+1)} - \rho_{XVW}^{(t+1)}\|_{\text{Tr}} \leq \mathbf{E}_{a \in_R \mathcal{A}} \|\tau_{XVW}^{(t,a)} - \rho_{XVW}^{(t)}\|_{\text{Tr}}$ .

*Proof.* Recall that

$$\begin{aligned} \rho_{XVW}^{(t+1)} &= \mathbf{E}_{a \in_R \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t,a,M(a,x)}(\rho_{VW|x}^{(t)}) \right], \\ \tau_{XVW}^{(t+1)} &= \mathbf{E}_{a \in_R \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t,a,M(a,x)}(\tau_{VW|x}^{(t,a)}) \right]. \end{aligned}$$

Therefore by triangle inequality and contractivity of quantum channels under trace norms,

$$\begin{aligned} \|\tau_{XVW}^{(t+1)} - \rho_{XVW}^{(t+1)}\|_{\text{Tr}} &\leq \mathbf{E}_{a \in_R \mathcal{A}} \left\| \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \left( \Phi_{t,a,M(a,x)}(\tau_{VW|x}^{(t,a)}) - \Phi_{t,a,M(a,x)}(\rho_{VW|x}^{(t)}) \right) \right\|_{\text{Tr}} \\ &= \mathbf{E}_{a \in_R \mathcal{A}} \sum_{x \in \mathcal{X}} \left\| \Phi_{t,a,M(a,x)}(\tau_{VW|x}^{(t,a)}) - \Phi_{t,a,M(a,x)}(\rho_{VW|x}^{(t)}) \right\|_{\text{Tr}} \\ &\leq \mathbf{E}_{a \in_R \mathcal{A}} \sum_{x \in \mathcal{X}} \|\tau_{VW|x}^{(t,a)} - \rho_{VW|x}^{(t)}\|_{\text{Tr}} \\ &= \mathbf{E}_{a \in_R \mathcal{A}} \|\tau_{XVW}^{(t,a)} - \rho_{XVW}^{(t)}\|_{\text{Tr}}. \quad \square \end{aligned}$$

#### 4.4 Proof of Theorem 2

*Proof.* First, combining Corollaries 4.6, 4.8, 4.10 and 4.13 and Lemma 4.14 we have

$$\|\tau_{XVW}^{(t+1)} - \rho_{XVW}^{(t+1)}\|_{\text{Tr}} \leq \|\tau_{XVW}^{(t)} - \rho_{XVW}^{(t)}\|_{\text{Tr}} + 8 \cdot 2^{q-2r} + 2^{-5r} + 2^{-2r}.$$

Since  $\tau_{XVW}^{(0)} = \rho_{XVW}^{(0)}$ , by triangle inequality we know that  $\|\tau_{XVW}^{(T)} - \rho_{XVW}^{(T)}\|_{\text{Tr}} \leq T \cdot 10 \cdot 2^{q-2r} \leq 10 \cdot 2^{q-r}$ , and thus

$$\|\tau_{XVW}^{(T,\infty)} - \rho_{XVW}^{(T)}\|_{\text{Tr}} \leq 10 \cdot 2^{q-r} + 8 \cdot 2^{q-2r} + 2^{-5r}.$$

This bounds the difference between the measurement probabilities of  $\tau_{XVW}^{(T,\infty)}$  and  $\rho_{XVW}^{(T)}$  under any measurement, specifically the difference between the success probability of the branching program  $\rho$  and the following value on  $\tau$ :

$$\sum_{\substack{x \in \mathcal{X}, v \in \{0,1\}^q, w \in \mathcal{W} \\ \tilde{x}(v,w)=x}} \langle x, v, w | \tau_{XVW}^{(T,\infty)} | x, v, w \rangle = \sum_{v \in \{0,1\}^q, w \in \mathcal{W}} \text{Tr}[\tau_{X|v,w}^{(T,\infty)}] \cdot P_{X|v,w}^{\tau_{X|v,w}^{(T,\infty)}}(\tilde{x}(v,w)).$$

Since  $\|P_{X|v,w}^{\tau_{X|v,w}^{(T,\infty)}}\|_{\infty} \leq 2^{2\ell+9r} \cdot 2^{-n}$  and  $\text{Tr}[\tau_{XVW}^{(T,\infty)}] \leq 1$ , the above value is at most  $2^{2\ell+9r} \cdot 2^{-n}$ . Therefore the success probability of the branching program  $\rho$  is at most (recall that  $2\ell + 9r - n \leq -r$ )

$$10 \cdot 2^{q-r} + 8 \cdot 2^{q-2r} + 2^{-5r} + 2^{2\ell+9r} \cdot 2^{-n} = O(2^{q-r}). \quad \square$$

## 5 Proof of Lemma 4.5

The first step towards proving Lemma 4.5 is to analyze how  $P_{X|v,w}^{\tau^{(t)}}$  evolves according to the rule

$$\tau_{XVW}^{(t+1)} = \mathbf{E}_{a \in \mathcal{R}^{\mathcal{A}}} \left[ \sum_{x \in \mathcal{X}} |x\rangle \langle x| \otimes \Phi_{t,a,M(a,x)}(\tau_{VW|x}^{(t,a)}) \right].$$

We introduce the following notations. For every  $a \in \mathcal{A}$  and  $b \in \{-1, 1\}$ , let

$$\mathbb{1}_{a,b} = \frac{1}{2}(\vec{1} + b \cdot M_a),$$

which is a 0-1 vector that indicates whether  $M(a, x) = b$ . Let

$$\tau_{XVW}^{(t,a,b)} = (\text{Diag } \mathbb{1}_{a,b} \otimes \mathbb{1}_{VW}) \tau_{XVW}^{(t,a)}, \quad (7)$$

so that we can write

$$\tau_{XVW}^{(t+1)} = \mathbf{E}_{a \in \mathcal{R}^{\mathcal{A}}} \left[ (\mathbb{1}_X \otimes \Phi_{t,a,1})(\tau_{XVW}^{(t,a,1)}) + (\mathbb{1}_X \otimes \Phi_{t,a,-1})(\tau_{XVW}^{(t,a,-1)}) \right]. \quad (8)$$

Thus Claim 3.3 implies that  $P_{X|v,w}^{\tau^{(t+1)}}$  is a convex combination of  $P_{X|v',w'}^{\tau^{(t,a,b)}}$  for some  $a, b, w'$  and  $|v'\rangle$ .

### 5.1 Target Distribution and Badness

Before considering the target distribution, let us first establish that the  $\ell_2$ -norms of  $P_{X|v,w}^{\tau^{(t)}}$  cannot be too large:

**Lemma 5.1.** *For every  $0 \leq t \leq T$ ,  $|v\rangle \in \mathcal{V}$ ,  $w \in \mathcal{W}$ , we have*

$$\|P_{X|v,w}^{\tau^{(t)}}\|_2 \leq 4 \cdot 2^\ell \cdot 2^{-n/2}.$$

*Proof.* When  $t = 0$  the statement is clearly true as  $P_{X|v,w}^{\tau^{(0)}}$  is always uniform.

Now assume  $t > 0$ . By Lemma 4.1 and Lemma 4.11 we know that

$$\|P_{X|v',w'}^{\tau^{(t-1,a)}}\|_2 \leq (1 - 2^{-r})^{-1} \cdot 2^\ell \cdot 2^{-n/2}$$

for every  $w' \in \mathcal{W}$ ,  $|v'\rangle \in \mathcal{V}$  and  $a \in \mathcal{A}$ , as  $\tau_{XVW}^{(t-1,a)}$  is truncated from  $\tau_{XVW}^{(t-1,\infty)}$ . Since  $G_a(P_{X|v',w'}^{\tau^{(t-1,a)}}$ ) is true, meaning that the distribution is evenly divided by  $a$ , we further have

$$\|P_{X|v',w'}^{\tau^{(t-1,a,b)}}\|_2 = \frac{\|\mathbb{1}_{a,b} \cdot P_{X|v',w'}^{\tau^{(t-1,a)}}\|_2}{\|\mathbb{1}_{a,b} \cdot P_{X|v',w'}^{\tau^{(t-1,a)}}\|_1} \leq 2(1 - 2^{-r})^{-1} \cdot \|P_{X|v',w'}^{\tau^{(t-1,a)}}\|_2 \leq 4 \cdot 2^\ell \cdot 2^{-n/2}.$$

Since  $P_{X|v,w}^{\tau^{(t)}}$  is a convex combination of  $P_{X|v',w'}^{\tau^{(t-1,a,b)}}$ , by convexity its  $\ell_2$ -norm is bounded by  $4 \cdot 2^\ell \cdot 2^{-n/2}$ .  $\square$



From now on we use  $P$  to denote a fixed target distribution (which we will later choose to be the distribution in Lemma 4.5), such that

$$2^\ell \cdot 2^{-n/2} \leq \|P\|_2 \leq 4 \cdot 2^\ell \cdot 2^{-n/2}.$$

We want to bound the progress of  $\langle P_{X|v,w}^{\tau^{(t)}}, P \rangle$ , which starts off as  $2^{-n}$  at  $t = 0$ , and becomes at least  $2^{2\ell} \cdot 2^{-n}$  when  $P_{X|v,w}^{\tau^{(t)}} = P$ . Note that by Cauchy-Schwarz we always have

$$\langle P_{X|v,w}^{\tau^{(t)}}, P \rangle \leq \|P_{X|v,w}^{\tau^{(t)}}\|_2 \|P\|_2 \leq 16 \cdot 2^{2\ell} \cdot 2^{-n}. \quad (9)$$

In order to bound the progress, we introduce some new notations. For any superscript (such as  $(t, a)$ ) on the partial systems, we use  $\sigma_{XVW}$  to denote  $\tau_{XVW}(\text{Diag } P \otimes \mathbb{I}_{VW})$ . Notice that

$$\text{Tr}[\sigma_{X|v,w}] = \text{Tr}[\tau_{X|v,w} \text{Diag } P] = \text{Tr}[\tau_{X|v,w}] \cdot \langle P_{X|v,w}^\tau, P \rangle.$$

Similarly,  $P_{X|v,w}^\sigma$  can be deduced from  $P_{X|v,w}^\tau$  via

$$P_{X|v,w}^\sigma(x) = \frac{\text{Tr}[\tau_{X|v,w}]}{\text{Tr}[\sigma_{X|v,w}]} \cdot P_{X|v,w}^\tau(x) \cdot P(x) = \frac{P_{X|v,w}^\tau(x) \cdot P(x)}{\langle P_{X|v,w}^\tau, P \rangle}. \quad (10)$$

Therefore we can bound the  $\ell_2$  norm of  $P_{X|v,w}^\sigma$  as

$$\|P_{X|v,w}^\sigma\|_2 \leq \frac{1}{\langle P_{X|v,w}^\tau, P \rangle} \cdot \|P_{X|v,w}^\tau\|_\infty \cdot \|P\|_2.$$

Now we can identify the places where  $\langle P_{X|v,w}^{\tau^{(t)}}, P \rangle$  increases by a lot, which happens when the *inner product* is not evenly divided by some  $a \in \mathcal{A}$  (we will see the reason in the analysis later). Formally, at stage  $0 \leq t < T$ , we say  $(w, a)$  is *bad* if

$$\exists |v\rangle \in \mathcal{V}, \text{ s.t. } |\langle M_a, P_{X|v,w}^{\sigma^{(t,a)}} \rangle| > 2^{-r} \text{ and } \langle P_{X|v,w}^{\tau^{(t,a)}}, P \rangle \geq \frac{1}{2} \cdot 2^{-n}. \quad (11)$$

**Lemma 5.2.** *For every  $0 \leq t < T$  and  $w \in \mathcal{W}$ , we have*

$$\Pr_{a \in_{\mathbb{R}} \mathcal{A}} [(w, a) \text{ is bad}] \leq 2^{-k}.$$

*Proof.* Since  $\tau_{XVW}^{(t,a)}$  is truncated from  $\tau_{XVW}^{(t,\infty)}$ , Lemma 4.1 shows that for every  $|v\rangle \in \mathcal{V}$ ,  $w \in \mathcal{W}$  and  $a \in \mathcal{A}$  there is  $|v'\rangle \in \mathcal{V}$  such that

$$P_{X|v,w}^{\tau^{(t,a)}} = P_{X|v',w}^{\tau^{(t,\infty)}}$$

and by Eq. (10) it also implies that

$$P_{X|v,w}^{\sigma^{(t,a)}} = P_{X|v',w}^{\sigma^{(t,\infty)}}.$$

Now fix some  $w \in \mathcal{W}$ , and let  $\mathcal{A}' \subseteq \mathcal{A}$  be the set of  $a \in \mathcal{A}$  such that

$$\exists |v\rangle \in \mathcal{V}, \text{ s.t. } |\langle M_a, P_{X|v,w}^{\sigma^{(t,\infty)}} \rangle| > 2^{-r} \text{ and } \langle P_{X|v,w}^{\tau^{(t,\infty)}}, P \rangle \geq \frac{1}{2} \cdot 2^{-n}.$$

Then  $\mathcal{A}'$  contains all  $a$  such that  $(w, a)$  is bad, and our goal is to bound the fraction of  $\mathcal{A}'$  in  $\mathcal{A}$ .

In the rest of the proof we temporarily omit the super script and write  $\tau^{(t,\infty)}$  and  $\sigma^{(t,\infty)}$  simply as  $\tau$  and  $\sigma$ . For the same reason as in Lemma 4.12 we can assume that  $\tau_{V|w} = \mathbb{I}_V$ , and thus

$$\langle v | \sigma_{V|w} | v \rangle = \text{Tr}[\sigma_{X|v,w}] = \langle P_{X|v,w}^\tau, P \rangle, \text{ and } \text{Tr}[\sigma_{XV|w}] = \langle P_{X|w}^\tau, P \rangle \leq 16 \cdot 2^{2\ell} \cdot 2^{-n}.$$

where the last inequality is by Lemma 4.11 and Cauchy-Schwarz, in the same way as Eq. (9).

Suppose that we have diagonalization  $\sigma_{V|w} = U^\dagger D U$ , where  $U$  is unitary and  $D$  is diagonal and non-negative. Let  $\mathcal{V}' \subseteq \mathcal{V}$  be the subspace spanned by  $U^\dagger |e\rangle$  over the computational basis vectors  $|e\rangle \in \mathcal{V}$  such that  $\langle e|D|e\rangle \geq 2^{-4r} \cdot 2^{-2\ell} \cdot 2^{-n}$ . So for every  $|v\rangle \in \mathcal{V}'$  we have

$$\langle P_{X|v,w}^\tau, P \rangle = \text{Tr}[\sigma_{X|v,w}] \geq 2^{-4r} \cdot 2^{-2\ell} \cdot 2^{-n}.$$

We claim that for every  $a \in \mathcal{A}'$ , there exists  $|v\rangle \in \mathcal{V}'$  such that  $|\langle M_a, P_{X|v,w}^\sigma \rangle| > \frac{1}{2} \cdot 2^{-r}$ . To prove the claim, let  $\Pi$  be the projection operator from  $\mathcal{V}$  to  $\mathcal{V}'$ , and then  $(\mathbb{I}_X \otimes \Pi)\sigma_{XV|w}(\mathbb{I}_X \otimes \Pi)$  can be conceptually seen as a truncated partial system  $\sigma_{XV|w}^{|G}$  where  $G(v, w)$  holds when  $\text{Tr}[\sigma_{X|v,w}] \geq 2^{-4r-2\ell} \cdot 2^{-n}$  for the fixed  $w$ . By Lemma 4.3 we have

$$\|\sigma_{XV|w}^{|G} - \sigma_{XV|w}\|_{\text{Tr}} \leq 3 \cdot 2^q \cdot \sqrt{2^{-4r-2\ell-n} \cdot \text{Tr}[\sigma_{XV|w}]} \leq 12 \cdot 2^{q-2r} \cdot 2^{-n}.$$

Since  $a \in \mathcal{A}'$ , assume for  $|u\rangle \in \mathcal{V}$  we have  $|\langle M_a, P_{X|u,w}^\sigma \rangle| > 2^{-r}$  and  $\text{Tr}[\sigma_{X|u,w}] = \langle P_{X|u,w}^\tau, P \rangle \geq \frac{1}{2} \cdot 2^{-n}$ . Let  $|v\rangle \sim \Pi|u\rangle$ , then we have

$$\begin{aligned} \|P_{X|u,w}^\sigma - P_{X|v,w}^\sigma\|_1 &= \|P_{X|u,w}^\sigma - P_{X|u,w}^{|G}\|_1 \leq \left\| \frac{\sigma_{X|u,w}}{\text{Tr}[\sigma_{X|u,w}]} - \frac{\sigma_{X|u,w}^{|G}}{\text{Tr}[\sigma_{X|u,w}^{|G}]}\right\|_{\text{Tr}} \\ &\leq \left\| \frac{\sigma_{X|u,w}}{\text{Tr}[\sigma_{X|u,w}]} - \frac{\sigma_{X|u,w}^{|G}}{\text{Tr}[\sigma_{X|u,w}]} \right\|_{\text{Tr}} + \left\| \frac{\sigma_{X|u,w}^{|G}}{\text{Tr}[\sigma_{X|u,w}]} - \frac{\sigma_{X|u,w}^{|G}}{\text{Tr}[\sigma_{X|u,w}^{|G}]}\right\|_{\text{Tr}} \\ &= \left\| \frac{\sigma_{X|u,w}}{\text{Tr}[\sigma_{X|u,w}]} - \frac{\sigma_{X|u,w}^{|G}}{\text{Tr}[\sigma_{X|u,w}]} \right\|_{\text{Tr}} + \left| \frac{1}{\text{Tr}[\sigma_{X|u,w}]} - \frac{1}{\text{Tr}[\sigma_{X|u,w}^{|G}]}\right| \cdot \text{Tr}[\sigma_{X|u,w}^{|G}|] \\ &= \frac{\|\sigma_{X|u,w} - \sigma_{X|u,w}^{|G}\|_{\text{Tr}}}{\text{Tr}[\sigma_{X|u,w}]} + \frac{|\text{Tr}[\sigma_{X|u,w}^{|G}] - \text{Tr}[\sigma_{X|u,w}]|}{\text{Tr}[\sigma_{X|u,w}]} \\ &\leq \frac{2\|\sigma_{X|u,w} - \sigma_{X|u,w}^{|G}\|_{\text{Tr}}}{\text{Tr}[\sigma_{X|u,w}]} \leq \frac{2\|\sigma_{XV|w} - \sigma_{XV|w}^{|G}\|_{\text{Tr}}}{\text{Tr}[\sigma_{X|u,w}]} \leq 48 \cdot 2^{q-2r} \leq \frac{1}{2} \cdot 2^{-r}, \end{aligned}$$

where the last step is due to  $q \leq r - 7$ . Thus

$$|\langle M_a, P_{X|v,w}^\sigma \rangle| \geq |\langle M_a, P_{X|u,w}^\sigma \rangle| - \|P_{X|u,w}^\sigma - P_{X|v,w}^\sigma\|_1 > \frac{1}{2} \cdot 2^{-r}.$$

Similarly to the proof for Lemma 4.12, for each  $a \in \mathcal{A}'$  let

$$\pi_a = \text{Tr}_X[(\text{Diag } M_a \otimes U^\dagger D^{-1/2} U) \cdot \sigma_{XV|w} \cdot (\mathbb{I}_X \otimes U^\dagger D^{-1/2} U)]$$

which is a Hermitian operator on  $\mathcal{V}$ . For each  $|v\rangle \in \mathcal{V}$ , let  $|v'\rangle \sim U^\dagger D^{-1/2} U|v\rangle$ . Recall that  $\sigma_{V|w} = U^\dagger D U$ , and therefore

$$\begin{aligned} P_{X|v,w}^\sigma &= \frac{\text{diag}(\mathbb{I}_X \otimes \langle v|)\sigma_{XV|w}(\mathbb{I}_X \otimes |v\rangle)}{\langle v|\sigma_{V|w}|v\rangle} \\ &= \frac{\text{diag}(\mathbb{I}_X \otimes \langle v'|U^\dagger D^{-1/2} U)\sigma_{XV|w}(\mathbb{I}_X \otimes U^\dagger D^{-1/2} U|v'\rangle)}{\langle v'|U^\dagger D^{-1/2} U\sigma_{V|w}U^\dagger D^{-1/2} U|v'\rangle} \\ &= \text{diag}(\mathbb{I}_X \otimes \langle v'|U^\dagger D^{-1/2} U)\sigma_{XV|w}(\mathbb{I}_X \otimes U^\dagger D^{-1/2} U|v'\rangle). \end{aligned}$$

And that means

$$\langle v' | \pi_a | v' \rangle = \left\langle M_a, \text{diag} (\mathbb{I}_X \otimes \langle v' | U^\dagger D^{-1/2} U \rangle \sigma_{XV|w} (\mathbb{I}_X \otimes U^\dagger D^{-1/2} U | v' \rangle) \right\rangle = \langle M_a, P_{X|v,w}^\sigma \rangle.$$

We showed above that there exists  $|v\rangle \in \mathcal{V}$ , and thus  $|v'\rangle \in \mathcal{V}'$  such that

$$|\langle v' | \pi_a | v' \rangle| = \left| \langle M_a, P_{X|v,w}^\sigma \rangle \right| \geq \frac{1}{2} \cdot 2^{-r},$$

which means that for  $\Pi \pi_a \Pi$ , the restriction of  $\pi_a$  on  $\mathcal{V}'$ , we have  $\|\Pi \pi_a \Pi\|_2 \geq \frac{1}{2} \cdot 2^{-r}$ . Now consider a uniformly random unit vector  $|v'\rangle$  in  $\mathcal{V}'$ , and by Lemma 3.1 we know that for some absolute constant  $c$ ,

$$\Pr_{|v'\rangle} \left[ |\langle v' | \sigma_a | v' \rangle| \geq 2^{-r'} \right] \geq 1 - 2^{(q+r+1-r')/2} c - e^{-2^q} \geq 1 - 2^{-r} c - e^{-1} \geq \frac{1}{2}.$$

Therefore, for the random vector  $|v\rangle \sim U^\dagger D^{-1/2} U |v'\rangle$  where  $|v'\rangle$  is uniform in  $\mathcal{V}'$ , we conclude that

$$\Pr_{|v\rangle} \left[ |\langle M_a, P_{X|v,w}^\sigma \rangle| \geq 2^{-r'} \right] \geq \frac{1}{2}.$$

On the other hand, as  $|v'\rangle \in \mathcal{V}'$ , it also holds that  $|v\rangle \in \mathcal{V}$ , therefore  $\langle P_{X|v,w}^\tau, P \rangle \geq 2^{-4r} \cdot 2^{-2\ell} \cdot 2^{-n}$  is always true. Thus there exists a  $|v\rangle \in \mathcal{V}$  that simultaneously satisfies

$$\langle P_{X|v,w}^\tau, P \rangle \geq 2^{-4r} \cdot 2^{-2\ell} \cdot 2^{-n} \quad \text{and} \quad |\langle M_a, P_{X|v,w}^\sigma \rangle| \geq 2^{-r'}$$

for at least  $1/2$  of  $a \in \mathcal{A}'$ . Since

$$\|P_{X|v,w}^\sigma\|_2 \leq \frac{1}{\langle P_{X|v,w}^\tau, P \rangle} \cdot \|P_{X|v,w}^\tau\|_\infty \cdot \|P\|_2 \leq 4 \cdot 2^{5\ell+13r} \cdot 2^{-n/2} = 2^{\ell'} \cdot 2^{-n/2},$$

and  $M$  is a  $(k', \ell')$ -extractor with error  $2^{-r'}$ , there are at most  $2^{-k'}$  fraction of  $a \in \mathcal{A}$  such that  $|\langle M_a, P_{X|v,w}^\sigma \rangle| \geq 2^{-r'}$ , which means that

$$\Pr_{a \in \mathcal{R}^{\mathcal{A}}} [(w, a) \text{ is bad}] \leq \Pr_{a \in \mathcal{R}^{\mathcal{A}}} [a \in \mathcal{A}'] \leq 2 \cdot 2^{-k'} = 2^{-k}. \quad \square$$

## 5.2 Badness Levels

At stage  $t$ , for each classical memory state  $w \in \mathcal{W}$  we count how many times the path to it has been bad, which is a random variable depending on the previous random choices of  $a \in \mathcal{A}$ . This is stored in another classical register  $B$ , which we call *badness level* and takes values  $\beta \in \{0, \dots, T\}$ . It is initially set to be 0, that is, we let

$$\tau_{XVWB}^{(0)} = \tau_{XVW}^{(0)} \otimes |0\rangle\langle 0|_B.$$

We ensure that the distribution of  $B$  always only depends on  $W$  and is independent of  $X$  and  $V$  conditioned on  $W$ , using the following updating rules on the combined system  $\tau_{XVWB}$  for each stage  $0 \leq t < T$ :

- The truncation steps are executed independently of  $B$ . Therefore, for each  $a \in \mathcal{A}$  we let

$$\tau_{XVWB}^{(t,a)} = \sum_{w \in \mathcal{W}} \tau_{XV|w}^{(t,a)} \otimes |w\rangle\langle w| \otimes \text{Diag } P_B^{T|w}. \quad (12)$$

- The value of  $B$  updates before the evolution step, where for each  $a \in \mathcal{A}$  and  $b \in \{-1, 1\}$  we let

$$\tau_{XVWB}^{(t,a,b)} = (\text{Diag } \mathbb{1}_{a,b} \otimes \mathbb{I}_V \otimes U_a) \tau_{XVWB}^{(t,a)} (\mathbb{I}_{XV} \otimes U_a^\dagger).$$

Here  $U_a$  is a permutation operator, depending on  $\tau_{XVW}^{(t,a)}$ , acting on  $\mathcal{W} \otimes \{0, \dots, T\}$  such that

$$U_a |w\rangle |\beta\rangle = \begin{cases} |w\rangle |(\beta + 1) \bmod (T + 1)\rangle & \text{if } (w, a) \text{ is bad,} \\ |w\rangle |\beta\rangle & \text{otherwise.} \end{cases}$$

- For the evolution step, we apply the channels  $\Phi_{t,a,b}$  on the memories  $W$  and  $V$  to get

$$\tau_{XVWB}^{(t+1)} = \mathbf{E}_{a \in \mathcal{R}\mathcal{A}} \left[ (\mathbb{I}_X \otimes \Phi_{t,a,1} \otimes \mathbb{I}_B) (\tau_{XVWB}^{(t,a,1)}) + (\mathbb{I}_X \otimes \Phi_{t,a,-1} \otimes \mathbb{I}_B) (\tau_{XVWB}^{(t,a,-1)}) \right].$$

Notice that the evolution step might introduce dependencies between  $X, V$  and  $B$ . However, such dependencies are eliminated later due to how we handle the truncation steps (12), and thus do not affect our proof.

We can check that the combined partial system  $\tau_{XVWB}^{(t)}$  defined above is consistent with the partial system  $\tau_{XVW}^{(t)}$  that we discussed in previous sections, in the sense that  $\text{Tr}_B[\tau_{XVWB}^{(t)}] = \tau_{XVW}^{(t)}$  always holds:

- For the truncation step, it is straightforward to check that

$$\text{Tr}_B[\tau_{XVWB}^{(t,a)}] = \sum_{w \in \mathcal{W}} \tau_{XV|w}^{(t,a)} \otimes |w\rangle \langle w| = \tau_{XVW}^{(t,a)}.$$

- The permutation operator  $U_a$  acts on  $\mathcal{W}$  as identity since

$$\text{Tr}_B \left[ U_a |w, \beta\rangle \langle w, \beta| U_a^\dagger \right] = |w\rangle \langle w|.$$

Recalling Eq. (7) that  $\tau_{XVW}^{(t,a,b)} = (\text{Diag } \mathbb{1}_{a,b} \otimes \mathbb{I}_V) \tau_{XVW}^{(t,a)}$ , we have  $\text{Tr}_B[\tau_{XVWB}^{(t,a,b)}] = \tau_{XVW}^{(t,a,b)}$ .

- The evolution step can be checked directly from the formula without  $B$  (Eq. (8)):

$$\tau_{XVW}^{(t+1)} = \mathbf{E}_{a \in \mathcal{R}\mathcal{A}} \left[ (\mathbb{I}_X \otimes \Phi_{t,a,1}) (\tau_{XVW}^{(t,a,1)}) + (\mathbb{I}_X \otimes \Phi_{t,a,-1}) (\tau_{XVW}^{(t,a,-1)}) \right].$$

So all previously proved properties about  $\tau_{XVW}^{(t)}$  are preserved. In addition, we prove the following two properties about badness levels.

**Lemma 5.3.** *For every  $0 \leq t \leq T$ ,  $|v\rangle \in \mathcal{V}$  and  $w \in \mathcal{W}$ , we have*

$$\langle P_{X|v,w}^{\tau^{(t)}} \rangle \leq \sum_{\beta=0}^T P_{B|w}^{\tau^{(t)}}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3t}.$$

*Proof.* We prove it by induction on  $t$ . For  $t = 0$  the lemma is true as  $\langle P_{X|v,w}^{\tau^{(0)}} \rangle = 2^{-n}$  and  $P_{B|w}^{\tau^{(0)}}(0) = 1$ .

Suppose the lemma holds true for some  $t < T$ . By a similar argument as in Lemma 4.11 and applying Lemma 4.1 multiple times, we know that for every  $|v\rangle \in \mathcal{V}$ ,  $w \in \mathcal{W}$  and  $a \in \mathcal{A}$ , there exists  $|v'\rangle$  and  $|v''\rangle \in \mathcal{V}$  such that

$$\langle P_{X|v,w}^{\tau^{(t,a)}} \rangle = \langle P_{X|v',w}^{\tau^{(t,\circ)}} \rangle \leq (1 - 2^{-r})^{-1} \langle P_{X|v',w}^{\tau^{(t,\star)}} \rangle = (1 - 2^{-r})^{-1} \langle P_{X|v'',w}^{\tau^{(t)}} \rangle,$$

and therefore

$$\langle P_{X|v,w}^{\tau(t,a)}, P \rangle \leq \sum_{\beta=0}^T P_{B|w}^{\tau(t)}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3t-1}. \quad (13)$$

Also, the truncation step by  $G_a$  implies that  $|\langle M_a, P_{X|v,w}^{\tau(t,a)} \rangle| \leq 2^{-r}$ . That is, for both  $b \in \{-1, 1\}$ ,

$$1 - 2^{-r} \leq 2 \|\mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}\|_1 \leq 1 + 2^{-r}.$$

Therefore we have, unconditionally

$$\langle P_{X|v,w}^{\tau(t,a,b)}, P \rangle = \frac{\langle \mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}, P \rangle}{\|\mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}\|_1} \leq 2(1 - 2^{-r})^{-1} \cdot \langle P_{X|v,w}^{\tau(t,a)}, P \rangle. \quad (14)$$

When the inner product is evenly divided, i.e.  $|\langle M_a, P_{X|v,w}^{\sigma(t,a)} \rangle| \leq 2^{-r}$ , we further have

$$\langle \mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}, P \rangle \leq \frac{1}{2}(1 + 2^{-r}) \langle P_{X|v,w}^{\tau(t,a)}, P \rangle \leq \frac{1}{2}(1 - 2^{-r})^{-1} \langle P_{X|v,w}^{\tau(t,a)}, P \rangle,$$

which means that

$$\langle P_{X|v,w}^{\tau(t,a,b)}, P \rangle = \frac{\langle \mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}, P \rangle}{\|\mathbb{1}_{a,b} \cdot P_{X|v,w}^{\tau(t,a)}\|_1} \leq (1 - 2^{-r})^{-2} \cdot \langle P_{X|v,w}^{\tau(t,a)}, P \rangle. \quad (15)$$

Now there are three cases to discuss:

- If  $(w, a)$  is bad, we have  $P_{B|w}^{\tau(t,a,b)}(\beta) = P_{B|w}^{\tau(t)}(\beta - 1)$  for every  $\beta > 0$ . Notice that  $P_{B|w}^{\tau(t)}(T) = 0$  as  $t < T$ , and thus Eq. (13) and Eq. (14) imply that

$$\begin{aligned} \langle P_{X|v,w}^{\tau(t,a,b)}, P \rangle &\leq \sum_{\beta=0}^{T-1} P_{B|w}^{\tau(t)}(\beta) \cdot 2^{\beta+1} \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3t-2} \\ &\leq \sum_{\beta=0}^T P_{B|w}^{\tau(t,a,b)}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3(t+1)}. \end{aligned}$$

- If  $(w, a)$  is not bad and  $|\langle M_a, P_{X|v,w}^{\sigma(t,a)} \rangle| \leq 2^{-r}$ , we have  $P_{B|w}^{\tau(t,a,b)}(\beta) = P_{B|w}^{\tau(t)}(\beta)$  for every  $\beta \geq 0$ . Then Eq. (13) and Eq. (15) imply that

$$\begin{aligned} \langle P_{X|v,w}^{\tau(t,a,b)}, P \rangle &\leq \sum_{\beta=0}^T P_{B|w}^{\tau(t)}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3t-3} \\ &= \sum_{\beta=0}^T P_{B|w}^{\tau(t,a,b)}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3(t+1)}. \end{aligned}$$

- If  $(w, a)$  is not bad and  $|\langle M_a, P_{X|v,w}^{\sigma(t,a)} \rangle| > 2^{-r}$ , by the definition of badness (11) we must have  $\langle P_{X|v,w}^{\tau(t,a)}, P \rangle < \frac{1}{2} \cdot 2^{-n}$ . Thus by Eq. (14),

$$\langle P_{X|v,w}^{\tau(t,a,b)}, P \rangle < (1 - 2^{-r})^{-1} \cdot 2^{-n} \leq \sum_{\beta=0}^T P_{B|w}^{\tau(t,a,b)}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3(t+1)}.$$

The last inequality follows from  $\sum_{\beta=0}^T P_{B|w}^{\tau^{(t,a,b)}}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3(t+1)} \geq 2^{-n}(1 - 2^{-r})^{-3(t+1)}$ . Hence we obtain the same conclusion from all three cases.

For the evolution step, since  $B$  is classical we can view  $X$  and  $B$  as a whole and apply Claim 3.3 on  $P_{XB|v,w}^{\tau^{(t+1)}}$ , which asserts that  $P_{XB|v,w}^{\tau^{(t+1)}}$  is a convex combination of  $P_{XB|v',w'}^{\tau^{(t,a,b)}}$  for some  $a, b, w'$  and  $|v'\rangle$ . Then by linearity we conclude that <sup>6</sup>

$$\langle P_{X|v,w}^{\tau^{(t+1)}}, P \rangle \leq \sum_{\beta=0}^T P_{B|w}^{\tau^{(t+1)}}(\beta) \cdot 2^\beta \cdot 2^{-n} \cdot (1 - 2^{-r})^{-3(t+1)}. \quad \square$$

**Lemma 5.4.** *For every  $0 \leq \beta \leq t \leq T$  we have*

$$\langle \beta | \tau_B^{(t)} | \beta \rangle \leq 2^{-k\beta} \binom{t}{\beta}.$$

*Proof.* We prove it by induction on  $t$ . For  $t = 0$  the lemma holds as  $\tau_B^{(0)} = |0\rangle\langle 0|_B$ . Also notice that the lemma is trivially true for every  $t$  when  $\beta = 0$ .

Now suppose the lemma holds for some  $t$ . By definition we have

$$\tau_B^{(t+1)} = \mathbf{E}_{a \in_R \mathcal{A}} [\tau_B^{(t,a,1)} + \tau_B^{(t,a,-1)}] = \mathbf{E}_{a \in_R \mathcal{A}} \text{Tr}_W [U_a \tau_{WB}^{(t,a)} U_a^\dagger].$$

Therefore

$$\langle \beta | \tau_B^{(t+1)} | \beta \rangle = \sum_{w \in \mathcal{W}} \mathbf{E}_{a \in_R \mathcal{A}} \left[ \langle w, \beta | U_a \tau_{WB}^{(t,a)} U_a^\dagger | w, \beta \rangle \right].$$

By Lemma 5.2 we know that for every  $w \in \mathcal{W}$ , the probability that  $(w, a)$  is bad for  $a \in_R \mathcal{A}$  is at most  $2^{-k}$ . In other words, for every  $\beta > 0$ ,

$$U_a^\dagger | w, \beta \rangle = \begin{cases} | w, \beta \rangle, & \text{w.p. } \geq 1 - 2^{-k} \\ | w, \beta - 1 \rangle, & \text{w.p. } \leq 2^{-k} \end{cases}$$

where the probability is taken over the random choice of  $a$ . It means that

$$\begin{aligned} \langle \beta | \tau_B^{(t+1)} | \beta \rangle &\leq \sum_{w \in \mathcal{W}} \langle w, \beta | \tau_{WB}^{(t,a)} | w, \beta \rangle + 2^{-k} \sum_{w \in \mathcal{W}} \langle w, \beta - 1 | \tau_{WB}^{(t,a)} | w, \beta - 1 \rangle \\ &= \langle \beta | \tau_B^{(t,a)} | \beta \rangle + 2^{-k} \cdot \langle \beta - 1 | \tau_B^{(t,a)} | \beta - 1 \rangle. \end{aligned}$$

Notice that

$$\tau_B^{(t,a)} = \sum_{w \in \mathcal{W}} \text{Tr}[\tau_{XV|w}^{(t,a)}] \cdot \text{Diag } P_{B|w}^{\tau^{(t)}} \leq \sum_{w \in \mathcal{W}} \text{Tr}[\tau_{XV|w}^{(t)}] \cdot \text{Diag } P_{B|w}^{\tau^{(t)}} = \tau_B^{(t)},$$

and thus we conclude that

$$\begin{aligned} \langle \beta | \tau_B^{(t+1)} | \beta \rangle &\leq \langle \beta | \tau_B^{(t)} | \beta \rangle + 2^{-k} \cdot \langle \beta - 1 | \tau_B^{(t)} | \beta - 1 \rangle \\ &\leq 2^{-k\beta} \binom{t}{\beta} + 2^{-k} \cdot 2^{-k(\beta-1)} \binom{t}{\beta-1} = 2^{-k\beta} \binom{t+1}{\beta}. \quad \square \end{aligned}$$

With the lemmas above in hand, we can finally prove Lemma 4.5.

<sup>6</sup>It should be noted that in  $\tau^{(t+1)}$ ,  $X$  and  $B$  are not independent. (In  $\tau^{(t,a,b)}$  they are independent (conditioned on  $v', w'$ )). Nevertheless, independence of  $X, B$  (in  $\tau^{(t+1)}$ ) is not needed or used here and we can conclude the final inequality by linearity by taking the corresponding convex combination of all inequalities.



*Proof for Lemma 4.5.* For the target distribution  $P = P_{X|v,w}^{\tau^{(t)}}$  we have  $\langle P_{X|v,w}^{\tau^{(t)}}, P \rangle > 2^{2\ell} \cdot 2^{-n}$ , so by Lemma 5.3,

$$\sum_{\beta=0}^T P_{B|w}^{\tau^{(t)}}(\beta) \cdot 2^\beta \cdot (1 - 2^{-r})^{-3t} > 2^{2\ell}.$$

Since  $t \leq T \leq 2^{r-2}$ , we have  $(1 - 2^{-r})^{-3t} \leq 2$ , and thus

$$\sum_{\beta=\ell}^T P_{B|w}^{\tau^{(t)}}(\beta) \cdot 2^\beta > \frac{1}{2} \left( 2^{2\ell} - 2 \cdot \sum_{\beta=0}^{\ell-1} 2^\beta \right) > 2^\ell.$$

On the other hand, for every  $\beta \geq \ell$ , by Lemma 5.4,

$$\text{Tr}[\tau_{B|w}^{(t)}] \cdot P_{B|w}^{\tau^{(t)}}(\beta) \leq \langle \beta | \tau_B^{(t)} | \beta \rangle \leq (2^{-kt})^\beta < 2^{-(k-r)\beta},$$

and thus by Eq. (6),

$$\text{Tr}[\tau_{X|v,w}^{(t)}] \leq \text{Tr}[\tau_{B|w}^{(t)}] < 2^{-\ell} \sum_{\beta=\ell}^T 2^{-(k-r)\beta} \cdot 2^\beta \leq 2 \cdot 2^{-(k-r)\ell} \leq 2^{-2m} \cdot 2^{-4r}. \quad \square$$

## Acknowledgement

We are grateful to Uma Girish for many important discussions and suggestions on the draft of this paper, and to the anonymous reviewers for their helpful comments.

## References

- [AAB<sup>+</sup>19] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019. 2
- [Aar18] Scott Aaronson. Shadow tomography of quantum states. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 325–338, 2018. 2
- [ACQ22] Dorit Aharonov, Jordan Cotler, and Xiao-Liang Qi. Quantum algorithmic measurement. *Nature communications*, 13(1):1–9, 2022. 2
- [ADR02] Yonatan Aumann, Yan Zong Ding, and Michael O Rabin. Everlasting security in the bounded storage model. *IEEE Transactions on Information Theory*, 48(6):1668–1680, 2002. 4
- [AR99] Yonatan Aumann and Michael O Rabin. Information theoretically secure communication in the limited storage space model. In *Annual International Cryptology Conference*, pages 65–79. Springer, 1999. 4
- [Aud07] Koenraad MR Audenaert. A sharp continuity estimate for the von neumann entropy. *Journal of Physics A: Mathematical and Theoretical*, 40(28):8127, 2007. 34
- [BCL20] Sebastien Bubeck, Sitan Chen, and Jerry Li. Entanglement is necessary for optimal quantum property testing. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 692–703. IEEE, 2020. 2

- [BGY18] Paul Beame, Shayan Oveis Gharan, and Xin Yang. Time-space tradeoffs for learning finite functions from random evaluations, with applications to polynomials. In *COLT*, volume 75 of *Proceedings of Machine Learning Research*, pages 843–856. PMLR, 2018. 2
- [BY21] Anne Broadbent and Peter Yuen. Device-independent oblivious transfer from the bounded-quantum-storage-model and computational assumptions. *arXiv preprint arXiv:2111.08595*, 2021. 4
- [CCHL22] Sitan Chen, Jordan Cotler, Hsin-Yuan Huang, and Jerry Li. Exponential separations between learning with and without quantum memory. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 574–585. IEEE, 2022. 2
- [CLO22] Sitan Chen, Jerry Li, and Ryan O’Donnell. Toward instance-optimal state certification with incoherent measurements. In *Conference on Learning Theory*, pages 2541–2596. PMLR, 2022. 2
- [CM97] Christian Cachin and Ueli Maurer. Unconditional security against memory-bounded adversaries. In *Annual International Cryptology Conference*, pages 292–306. Springer, 1997. 4
- [CW01] Anthony Carbery and James Wright. Distributional and  $L^q$  norm inequalities for polynomials over convex bodies in  $\mathbb{R}^n$ . *Mathematical Research Letters*, 8(3):233–248, 2001. 10
- [DFR<sup>+</sup>07] Ivan B Damgård, Serge Fehr, Renato Renner, Louis Salvail, and Christian Schaffner. A tight high-order entropic quantum uncertainty relation with applications. In *Annual International Cryptology Conference*, pages 360–378. Springer, 2007. 4
- [DFSS07] Ivan B Damgård, Serge Fehr, Louis Salvail, and Christian Schaffner. Secure identification and qkd in the bounded-quantum-storage model. In *Annual International Cryptology Conference*, pages 342–359. Springer, 2007. 4
- [DFSS08] Ivan B Damgård, Serge Fehr, Louis Salvail, and Christian Schaffner. Cryptography in the bounded-quantum-storage model. *SIAM Journal on Computing*, 37(6):1865–1890, 2008. 4
- [DM02] Stefan Dziembowski and Ueli Maurer. Tight security proofs for the bounded-storage model. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 341–350, 2002. 4
- [DM04] Stefan Dziembowski and Ueli Maurer. On generating the initial key in the bounded-storage model. In *International Conference on the Theory and Applications of Cryptographic Techniques*, pages 126–137. Springer, 2004. 4
- [DQW21] Yevgeniy Dodis, Willy Quach, and Daniel Wichs. Speak much, remember little: Cryptography in the bounded storage model, revisited. *Cryptology ePrint Archive*, 2021. 4
- [DQW22] Yevgeniy Dodis, Willy Quach, and Daniel Wichs. Authentication in the bounded storage model. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 737–766. Springer, 2022. 4

- [DR02] Yan Zong Ding and Michael O Rabin. Hyper-encryption and everlasting security. In *Annual Symposium on Theoretical Aspects of Computer Science*, pages 1–26. Springer, 2002. 4
- [FvdG99] Christopher A Fuchs and Jeroen van de Graaf. Cryptographic distinguishability measures for quantum-mechanical states. *IEEE Transactions on Information Theory*, 45(4):1216–1227, 1999. 31
- [GKK<sup>+</sup>08] Dmitry Gavinsky, Julia Kempe, Iordanis Kerenidis, Ran Raz, and Ronald de Wolf. Exponential separation for one-way quantum communication complexity, with applications to cryptography. *SIAM J. Comput.*, 38(5):1695–1708, 2008. 4
- [GKLR21] Sumegha Garg, Pravesh K. Kothari, Pengda Liu, and Ran Raz. Memory-sample lower bounds for learning parity with noise. In *APPROX-RANDOM*, volume 207 of *LIPICs*, pages 60:1–60:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021. 2
- [GKR20] Sumegha Garg, Pravesh K. Kothari, and Ran Raz. Time-space tradeoffs for distinguishing distributions and applications to security of goldreich’s PRG. In *APPROX-RANDOM*, volume 176 of *LIPICs*, pages 21:1–21:18. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020. 2
- [GRT18] Sumegha Garg, Ran Raz, and Avishay Tal. Extractor-based time-space lower bounds for learning. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2018, page 990–1002, New York, NY, USA, 2018. Association for Computing Machinery. 2, 3, 4, 6
- [GZ19] Jiaxin Guan and Mark Zhandary. Simple schemes in the bounded storage model. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 500–524. Springer, 2019. 4
- [HHJ<sup>+</sup>16] Jeongwan Haah, Aram W Harrow, Zhengfeng Ji, Xiaodi Wu, and Nengkun Yu. Sample-optimal tomography of quantum states. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 913–925, 2016. 2
- [HKP20] Hsin-Yuan Huang, Richard Kueng, and John Preskill. Predicting many properties of a quantum system from very few measurements. *Nature Physics*, 16(10):1050–1057, 2020. 2
- [HN06] Danny Harnik and Moni Naor. On everlasting security in the hybrid bounded storage model. In *International Colloquium on Automata, Languages, and Programming*, pages 192–203. Springer, 2006. 4
- [KK12] Roy Kasher and Julia Kempe. Two-source extractors secure against quantum adversaries. *Theory of Computing*, 8(1):461–486, 2012. 33
- [KRT17] Gillat Kol, Ran Raz, and Avishay Tal. Time-space hardness of learning sparse parities. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1067–1080, 2017. 2
- [LM00] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5), October 2000. 10

- [Lu02] Chi-Jen Lu. Hyper-encryption against space-bounded adversaries from on-line strong extractors. In *Annual International Cryptology Conference*, pages 257–271. Springer, 2002. 4
- [LV21] Jiahui Liu and Satyanarayana Vusirikala. Secure multiparty computation in the bounded storage model. In *IMA International Conference on Cryptography and Coding*, pages 289–325. Springer, 2021. 4
- [Mau92] Ueli M Maurer. Conditionally-perfect secrecy and a provably-secure randomized cipher. *Journal of Cryptology*, 5(1):53–66, 1992. 4
- [MM18] Dana Moshkovitz and Michal Moshkovitz. Entropy samplers and strong generic lower bounds for space bounded learning. In *ITCS*, volume 94 of *LIPICs*, pages 28:1–28:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018. 2
- [MSTS04] Tal Moran, Ronen Shaltiel, and Amnon Ta-Shma. Non-interactive timestamping in the bounded storage model. In *Annual International Cryptology Conference*, pages 460–476. Springer, 2004. 4
- [OP04] Masanori Ohya and Dénes Petz. *Quantum entropy and its use*. Springer Science & Business Media, 2004. 34
- [PMLA13] Stefano Pironio, Ll Masanes, Anthony Leverrier, and Antonio Acín. Security of device-independent quantum key distribution in the bounded-quantum-storage model. *Physical Review X*, 3(3):031007, 2013. 4
- [Raz17] Ran Raz. A time-space lower bound for a large class of learning problems. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 732–742, 2017. 2, 3, 4, 6
- [Raz18] Ran Raz. Fast learning requires good memory: A time-space lower bound for parity learning. *J. ACM*, 66(1), dec 2018. 2, 4
- [RS08] Alvin C Rencher and G Bruce Schaalje. *Linear models in statistics*. John Wiley & Sons, 2008. 10
- [Sch07] Christian Schaffner. Cryptography in the bounded-quantum-storage model. *arXiv preprint arXiv:0709.0289*, 2007. 4
- [Sha14] Ohad Shamir. Fundamental limits of online and distributed algorithms for statistical learning and estimation. *Advances in Neural Information Processing Systems*, 27, 2014. 2
- [SSV19] Vatsal Sharan, Aaron Sidford, and Gregory Valiant. Memory-sample tradeoffs for linear regression with small error. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 890–901, 2019. 2
- [SVW16] Jacob Steinhardt, Gregory Valiant, and Stefan Wager. Memory, communication, and statistical queries. In *Conference on Learning Theory*, pages 1490–1516. PMLR, 2016. 2
- [Uhl76] Armin Uhlmann. The “transition probability” in the state space of a  $*$ -algebra. *Reports on Mathematical Physics*, 9(2):273–279, 1976. 32

- [Wri16] John Wright. *How to learn a quantum state*. PhD thesis, Carnegie Mellon University, 2016. 2
- [WW08] Stephanie Wehner and Jürg Wullschleger. Composable security in the bounded-quantum-storage model. In *International Colloquium on Automata, Languages, and Programming*, pages 604–615. Springer, 2008. 4

## A Bounding Parameters for Theorem 2

*Proof.* For Equation (4), since  $r \leq r'/4$  by the definition of  $r$  and  $q \leq r - 7$ , we have:

$$q + r + 1 - r' \leq q + r + 1 - 4r \leq -2r - 6 \leq -2r.$$

For Equation (5), we use the fact that  $\ell = \frac{1}{5}(\ell' - 13r - 2)$ , therefore

$$\begin{aligned} 2\ell + 9r - n &= \frac{2}{5}(\ell' - 13r - 2) + 9r - n \\ &= \frac{2}{5}\ell' + \frac{19}{5}r - \frac{4}{5} - n. \end{aligned}$$

Solving the inequality  $\frac{2}{5}\ell' + \frac{19}{5}r - \frac{4}{5} - n \leq -r$  (together with the fact that  $\ell' \leq n$ ) gives us  $r \leq \frac{1}{8}\ell' + \frac{1}{6}$ , which follows since  $r \leq \frac{1}{26}\ell' + \frac{1}{6}$  by definition.

For Equation (6), as  $r \leq \frac{1}{26}\ell' + \frac{1}{6}$ ,  $r \leq \frac{1}{2}(k' - 1)$  and  $k = k' - 1$ ,

$$\begin{aligned} k - r &\geq k - (k' - 1)/2 = (k' - 1)/2, \\ \ell &= \frac{1}{5}(\ell' - 13r - 2) \geq \frac{1}{5}\left(\ell' - \frac{\ell'}{2} - \frac{13}{6} - 2\right) > \frac{1}{10}\ell' - 1. \end{aligned}$$

When  $\ell'$  is sufficiently large, we have  $\frac{1}{10}\ell' - 1 > \frac{1}{11}\ell' + 5$ . Thereby, using the fact that  $m \leq (k' - 1)\ell'/44$ :

$$\begin{aligned} (k - r)\ell &\geq \frac{1}{2}(k' - 1)\left(\frac{1}{11}\ell' + 5\right) \\ &= \frac{1}{22}(k' - 1)\ell' + \frac{5}{2}(k' - 1) \\ &\geq 2m + 5r \\ &\geq 2m + 4r + 1. \end{aligned} \quad \square$$

## B Proof for Proposition 4.2

We first state a variant of Fuchs-van de Graaf inequality [FvdG99] on fidelity, defined for two partial density operators  $\rho$  and  $\sigma$  as

$$F(\rho, \sigma) = \text{Tr} \left[ \sqrt{\sqrt{\rho}\sigma\sqrt{\rho}} \right].$$

**Lemma B.1.** *Let  $\rho, \sigma$  be two semi-definite operators. Assume  $\text{Tr}[\rho] \geq \text{Tr}[\sigma]$ . Then*

$$\frac{1}{2}\|\rho - \sigma\|_{\text{Tr}} \leq \sqrt{\frac{1}{4}(\text{Tr}[\rho] + \text{Tr}[\sigma])^2 - F(\rho, \sigma)} \leq \sqrt{\text{Tr}[\rho]^2 - F(\rho, \sigma)}.$$

Notice that when  $\text{Tr}[\rho] = \text{Tr}[\sigma] = 1$ , the above inequality is the original Fuchs-van de Graaf inequality.

*Proof.* Let  $u$  and  $v$  be purifications of  $\rho$  and  $\sigma$ , that is,  $u, v \in \mathcal{H} \otimes \mathcal{H}_A$  with  $\rho = \text{Tr}_A[uu^\dagger]$  and  $\sigma = \text{Tr}_A[vv^\dagger]$ , where  $\mathcal{H}$  is the ambient space of  $\rho$  and  $\sigma$ , and  $\mathcal{H}_A$  is some finite-dimensional Hilbert space.

Let  $U$  be a unitary that diagonalizes  $uu^\dagger - vv^\dagger$ , that is there is a diagonal matrix  $\Lambda \in \mathbb{C}^{d \times d}$  such that  $uu^\dagger - vv^\dagger = U\Lambda U^\dagger$ . Let  $p, q \in \mathbb{R}_{\geq 0}^d$  be the diagonal of  $U^\dagger uu^\dagger U$  and  $U^\dagger vv^\dagger U$  respectively. We have

$$\begin{aligned} u^\dagger u &= \text{Tr}[U^\dagger uu^\dagger U] = \sum_{x \in [d]} p(x), \\ v^\dagger v &= \text{Tr}[U^\dagger vv^\dagger U] = \sum_{x \in [d]} q(x), \\ \|uu^\dagger - vv^\dagger\|_{\text{Tr}} &= \|\Lambda\|_{\text{Tr}} = \sum_{x \in [d]} |p(x) - q(x)|, \\ |\langle u, v \rangle| &= |\langle U^\dagger u, U^\dagger v \rangle| \leq \sum_{x \in [d]} \sqrt{p(x)q(x)}. \end{aligned}$$

Therefore, by Cauchy-Schwarz inequality,

$$\begin{aligned} \|uu^\dagger - vv^\dagger\|_{\text{Tr}}^2 &= \left( \sum_{x \in [d]} |p(x) - q(x)| \right)^2 \\ &= \left( \sum_{x \in [d]} \left| \sqrt{p(x)} - \sqrt{q(x)} \right| \cdot \left| \sqrt{p(x)} + \sqrt{q(x)} \right| \right)^2 \\ &\leq \left( \sum_{x \in [d]} \left| \sqrt{p(x)} - \sqrt{q(x)} \right|^2 \right) \left( \sum_{x \in [d]} \left| \sqrt{p(x)} + \sqrt{q(x)} \right|^2 \right) \\ &= \left( \sum_{x \in [d]} p(x) + \sum_{x \in [d]} q(x) \right)^2 - 4 \left( \sum_{x \in [d]} \sqrt{p(x)q(x)} \right)^2 \\ &\leq \left( u^\dagger u + v^\dagger v \right)^2 - 4|\langle u, v \rangle|^2. \end{aligned}$$

Notice that  $\|\rho - \sigma\|_{\text{Tr}} \leq \|uu^\dagger - vv^\dagger\|_{\text{Tr}}$ ,  $\text{Tr}[\rho] = u^\dagger u$  and  $\text{Tr}[\sigma] = v^\dagger v$ . By Uhlmann's theorem [Uhl76], we can also choose  $u$  and  $v$  such that  $F(\rho, \sigma) = |\langle u, v \rangle|^2$ . Plugging them into the above inequality concludes the proof.  $\square$

Now we are ready to prove Proposition 4.2

*Proof for Proposition 4.2.* By Fuchs-van de Graaf inequality, it suffices to prove the following bound on fidelity:

$$F(\rho, \Pi\rho\Pi) \geq \text{Tr}[\Pi\rho]^2.$$

Let  $u$  be a purification of  $\rho_{XV}$ , that is,  $\rho = \text{Tr}_A[uu^\dagger]$  for some Hilbert space  $\mathcal{H}_A$ . Then  $(\Pi \otimes \mathbb{I}_A)u$  is a purification of  $\Pi\rho\Pi$ . By Uhlmann's theorem we have

$$F(\rho, \Pi\rho\Pi) \geq \left| u^\dagger (\Pi \otimes \mathbb{I}_A) u \right|^2 = \text{Tr} \left[ (\Pi \otimes \mathbb{I}_A) uu^\dagger \right]^2 = \text{Tr} \left[ \Pi \cdot \text{Tr}_A[uu^\dagger] \right]^2 = \text{Tr}[\Pi\rho]^2. \quad \square$$

## C Linear Quantum Memory Lower Bound

In this appendix, we prove Theorem 3 that shows a simpler proof for a linear quantum-memory lower bound (without classical memory). While Theorem 3 is qualitatively weaker than our main results in most cases, as it only gives a lower bound for programs with a linear-size quantum memory but without a (possibly quadratic) classical memory, Theorem 3 is technically incomparable to the main results, as it's stated in terms of quantum extractors and the bound on the quantum-memory size depends on a different parameter of the extractor. Additionally, the proof of Theorem 3 is significantly simpler than the proof of our main theorem.

We first define the quantum extractor property that we need, which is a simplified version of the ones considered in [KK12]. Given a matrix  $M : \mathcal{A} \times \mathcal{X} \rightarrow \{-1, 1\}$ , consider two independent sources  $A$  and  $X$  uniformly distributed over  $\mathcal{A}$  and  $\mathcal{X}$  respectively. Suppose there is some quantum register  $V$  whose state depends on  $A$  and  $X$ , and they together form a classical-quantum system

$$\rho_{AXV} = \bigoplus_{a \in \mathcal{A}, x \in \mathcal{X}} \rho_{V|a,x},$$

where  $\rho_{V|a,x}$  is the state of  $V$  when  $A = a$  and  $X = x$ . For any function  $f$  on  $A \times X$ , we say that  $V$  depends only on  $f(A, X)$  if for any  $a, a' \in \mathcal{A}$  and  $x, x' \in \mathcal{X}$ , whenever  $f(a, x) = f(a', x')$  we have  $\rho_{V|a,x} = \rho_{V|a',x'}$ . In particular,  $V$  depending only on  $A$  is equivalent to  $V$  being independent of  $X$ , or  $\rho_{XV} = \rho_X \otimes \rho_V$ .

We say that  $M$  is an  $X$ -strong  $(q, r)$ -quantum extractor, if for every classical-quantum system  $\rho_{AXV}$ , as above, with the  $q$ -qubit quantum subsystem  $V$  that depends only on  $A$ , it holds that

$$\left\| \rho_{M(A,X)XV} - U \otimes \rho_X \otimes \rho_V \right\|_{\text{Tr}} \leq 2^{-r}.$$

Here  $U = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}$  is the uniform operator over one bit, and  $\rho_{M(A,X)XV}$  is the classical-quantum system constructed by adding a new classical register storing the value of  $M(A, X)$  and tracing out  $A$ . In other words,

$$\rho_{M(A,X)XV} = \bigoplus_{b \in \{-1, 1\}, x \in \mathcal{X}} \sum_{\substack{a \in \mathcal{A} \\ M(a,x)=b}} \rho_{V|a,x}.$$

Notice that if we choose  $V$  to be trivial, the above inequality immediately implies that  $|\mathbf{E}[M(A, X)]| \leq 2^{-r}$ .

As an example, the results in [KK12] imply that the inner product function on  $n$  bits, where  $\mathcal{A} = \mathcal{X} = \mathbb{F}_2^n$  and

$$M(a, x) = (-1)^{a \cdot x},$$

is an  $X$ -strong  $(k, n - k)$ -quantum extractor for every  $2 \leq k \leq n$ .

In this section we prove the following theorem:

**Theorem 3.** *Let  $\mathcal{X}, \mathcal{A}$  be two finite sets with  $n = \log_2 |\mathcal{X}|$ . Let  $M : \mathcal{A} \times \mathcal{X} \rightarrow \{-1, 1\}$  be a matrix which is a  $X$ -strong  $(q, r)$ -quantum extractor. Let  $\rho$  be a branching program for the learning problem corresponding to  $M$ , described by classical-quantum systems  $\rho_{XV}^{(t)}$ , with  $q/2$ -qubit quantum memory  $V$  and length  $T$ , and without classical memory. Then the success probability of  $\rho$  is at most*

$$2^{-n} + 8T\sqrt{n+q} \cdot 2^{-r/4}.$$

To prove the theorem, we first need to define the following measure of dependency:

**Definition.** Let  $\rho_{XV}$  be a classical-quantum system over classical  $X$  and quantum  $V$ . The dependency of  $V$  on  $X$  in  $\rho_{XV}$  is defined as

$$\xi^\rho(X; V) = \min_{\tau_V} \|\rho_{XV} - \rho_X \otimes \tau_V\|_{\text{Tr}}$$

where  $\tau_V$  is taken over all density operators on  $V$ . Notice that in this definition taking  $\tau_V = \rho_V$  is almost optimal as we have

$$\|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}} \leq \|\rho_{XV} - \rho_X \otimes \tau_V\|_{\text{Tr}} + \|\rho_V - \tau_V\|_{\text{Tr}} \leq 2\|\rho_{XV} - \rho_X \otimes \tau_V\|_{\text{Tr}}. \quad (16)$$

We also consider the quantum mutual information between  $X$  and  $V$  in  $\rho_{XV}$ , defined as

$$\mathbf{I}_\rho(X; V) = \mathbf{S}(\rho_X) + \mathbf{S}(\rho_V) - \mathbf{S}(\rho_{XV}) = \mathbf{S}(\rho_{XV} \parallel \rho_X \otimes \rho_V),$$

where  $\mathbf{S}(\cdot)$  denotes the von Neumann entropy, and  $\mathbf{S}(\cdot \parallel \cdot)$  denotes quantum relative entropy. When  $V$  consists of  $q$  qubits, we have the following relationship between our dependency measure and quantum mutual information:

**Lemma C.1.**  $\frac{1}{2}\xi^\rho(X; V)^2 \leq \mathbf{I}_\rho(X; V) \leq q \cdot \xi^\rho(X; V) + 2\sqrt{\xi^\rho(X; V)}$ .

*Proof.* On one hand, using the inequality on quantum relative entropy and trace distance (see e.g. [OP04, Theorem 1.15]), we have

$$\mathbf{I}_\rho(X; V) = \mathbf{S}(\rho_{XV} \parallel \rho_X \otimes \rho_V) \geq \frac{1}{2}\|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}}^2 \geq \frac{1}{2}\xi^\rho(X; V)^2.$$

On the other hand, Fannes-Audenaert inequality [Aud07] tells us that for every  $x \in \mathcal{X}$ , the difference between the von Neumann entropies of any two states  $\rho$  and  $\tau$  on  $V$  is bounded by

$$|\mathbf{S}(\rho) - \mathbf{S}(\tau)| \leq q \cdot \frac{1}{2}\|\rho - \tau\|_{\text{Tr}} + h\left(\frac{1}{2}\|\rho - \tau\|_{\text{Tr}}\right)$$

where  $h(\epsilon) = -\epsilon \log_2 \epsilon - (1 - \epsilon) \log_2(1 - \epsilon)$  is the binary entropy function. Since the state of  $V$  conditioned on  $X = x$  is  $\rho_{V|x} / \Pr[X = x] = 2^n \rho_{V|x}$ , we have

$$\begin{aligned} \mathbf{I}_\rho(X; V) &= \mathbf{E}_{x \sim X} \left[ \mathbf{S}(\rho_V) - \mathbf{S}(2^n \rho_{V|x}) \right] \\ &\leq \frac{1}{2}q \cdot \mathbf{E}_{x \sim X} \|\rho_V - 2^n \rho_{V|x}\|_{\text{Tr}} + \mathbf{E}_{x \sim X} h\left(\frac{1}{2}\|\rho_V - 2^n \rho_{V|x}\|_{\text{Tr}}\right) \\ &\leq \frac{1}{2}q \cdot \|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}} + h\left(\frac{1}{2}\|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}}\right) \\ &\leq \frac{1}{2}q \cdot \|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}} + \sqrt{2\|\rho_{XV} - \rho_X \otimes \rho_V\|_{\text{Tr}}}, \end{aligned}$$

as  $h$  is concave and  $h(\epsilon) \leq 2\sqrt{\epsilon}$ . Now let  $\tau_V$  be the optimal density operator in the definition of  $\xi^\rho(X; V)$ . Plugging in Eq. (16), we conclude that

$$\mathbf{I}_\rho(X; V) \leq q \cdot \xi^\rho(X; V) + 2\sqrt{\xi^\rho(X; V)}. \quad \square$$

**Lemma C.2.** For every classical-quantum system  $\rho_{AXV}$  with the  $q$ -qubit quantum subsystem  $V$  that depends only on  $A$ , we have

$$\mathbf{I}_\rho(X; M(A, X), V) \leq 2(n + q) \cdot 2^{-r/2}.$$



*Proof.* Since  $\mathbf{I}_\rho(X; V) = 0$ , it suffices to bound  $\mathbf{I}_\rho(X; M(A, X) | V) \leq \mathbf{I}_\rho(M(A, X); X, V)$ . To bound the later, we first notice that since  $M$  is a strong  $(q, r)$ -quantum extractor,

$$\begin{aligned} \xi^\rho(M(A, X); X, V) &\leq \left\| \rho_{M(A, X)XV} - \rho_{M(A, X)} \otimes \rho_X \otimes \rho_V \right\|_{\text{Tr}} \\ &\leq \left\| \rho_{M(A, X)XV} - U \otimes \rho_X \otimes \rho_V \right\|_{\text{Tr}} + |\mathbf{E}[M(A, X)]| \\ &\leq 2 \cdot 2^{-r}. \end{aligned}$$

As the total dimension of  $X$  and  $V$  is  $2^{n+q}$ , by Lemma C.1 we have

$$\begin{aligned} \mathbf{I}_\rho(X; M(A, X), V) &\leq \mathbf{I}_\rho(M(A, X); X, V) \\ &\leq (n+q) \cdot \xi^\rho(M(A, X); X, V) + 2\sqrt{\xi^\rho(M(A, X); X, V)} \\ &\leq 5(n+q) \cdot 2^{-r/2}. \end{aligned} \quad \square$$

**Lemma C.3.** *For every classical-quantum system  $\rho_{AXV}$  with  $q/2$ -qubit quantum subsystem  $V$  that depends only on  $A$  and  $M(A, X)$ , we have*

$$\xi^\rho(X; V) \leq 4\sqrt{n+q} \cdot 2^{-r/4}.$$

*Proof.* Let  $W = \rho_{a,0} \otimes \rho_{a,1}$ , where  $\rho_{a,b}$  is the density matrix of  $V$  when  $A = a$  and  $M(A, X) = b$ . Then  $W$  is a  $q$ -bit quantum system that depends only on  $A$ . Since  $V$  can be decided from  $M(A, X)$  and  $W$ , we have

$$\xi^\rho(X; V)^2 \leq 2\mathbf{I}_\rho(X; V) \leq 2\mathbf{I}_\rho(X; M(A, X), W) \leq 10(n+q) \cdot 2^{-r/2}. \quad \square$$

We are now ready to prove Theorem 3. Let  $\Phi_{t,a,b}$  be the quantum channel applied on  $V$  at stage  $t$  with sample  $(a, b)$ , and recall that the evolution of the system  $\rho_{XV}^{(t)}$  can be expressed as

$$\rho_{XV}^{(t+1)} = \mathbf{E}_{a \sim A} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t,a,M(a,x)}(\rho_{V|x}^{(t)}) \right].$$

*Proof of Theorem 3.* We are going to bound the increment of  $\xi_t$ , which is the shorthand for  $\xi^{\rho^{(t)}}(X; V)$ . For now let us focus on some stage  $t$ , and let  $\tau$  be the density operator that minimizes  $\xi_t = \left\| \rho_{XV}^{(t)} - \rho_X^{(t)} \otimes \tau \right\|_{\text{Tr}}$ . Notice that  $\rho_X^{(t)} = \rho_X = 2^{-n}\mathbb{I}_X$  for every  $t$ .

Since  $\tau$  is a fixed quantum state, we can prepare  $\tau$  and apply  $\Phi_{A,M(A,X)}$  on  $\tau$  to obtain a new quantum register  $V'$ , which depends only on  $A$  and  $M(A, X)$ . Notice that

$$\rho_{XV'} = \mathbf{E}_{a \sim A} \left[ \sum_{x \in \mathcal{X}} |x\rangle\langle x| \otimes \Phi_{t,a,M(a,x)}(\tau) \right],$$

and therefore similarly to Lemma 4.14 (that evolution does not increase trace distance), we can show that

$$\begin{aligned} \left\| \rho_{XV}^{(t+1)} - \rho_{XV'} \right\|_{\text{Tr}} &\leq \mathbf{E}_{a \sim A} \sum_{x \in \mathcal{X}} \left\| \Phi_{t,a,M(a,x)}(\rho_{V|x}^{(t)}) - \Phi_{t,a,M(a,x)}(\tau) \right\|_{\text{Tr}} \\ &\leq \sum_{x \in \mathcal{X}} \left\| \rho_{V|x}^{(t)} - \tau \right\|_{\text{Tr}} \leq \left\| \rho_{XV}^{(t)} - \rho_X \otimes \tau \right\|_{\text{Tr}} = \xi_t. \end{aligned}$$

Hence we have

$$\begin{aligned}
\xi_{t+1} &\leq \left\| \rho_{XV}^{(t+1)} - \rho_X \otimes \rho_{V'} \right\|_{\text{Tr}} \\
&\leq \left\| \rho_{XV}^{(t+1)} - \rho_{XV'} \right\|_{\text{Tr}} + \left\| \rho_{XV'} - \rho_X \otimes \rho_{V'} \right\|_{\text{Tr}} \\
&\leq \xi_t + 2\xi^\rho(X; V') \\
&\leq \xi_t + 8\sqrt{n+q} \cdot 2^{-r/4}.
\end{aligned}$$

Since  $\xi_0 = 0$ , we conclude that

$$\xi_T \leq 8T\sqrt{n+q} \cdot 2^{-r/4}.$$

This value bounds the difference of the success probability of  $\rho$ , and that of a quantum branching program whose memory is independent of  $X$ . The later is clearly at most  $2^{-n}$ , which finishes the proof.  $\square$