



Cumulative Memory Lower Bounds for Randomized and Quantum Computation

Paul Beame   

Computer Science & Engineering, University of Washington, Seattle, WA, USA

Niels Kornerup   

Computer Science, University of Texas, Austin, TX, USA

Abstract

Cumulative memory – the sum of space used per step over the duration of a computation – is a fine-grained measure of time-space complexity that was introduced to analyze cryptographic applications like password hashing. It is a more accurate cost measure for algorithms that have infrequent spikes in memory usage and are run in environments such as cloud computing that allow dynamic allocation and de-allocation of resources during execution, or when many multiple instances of an algorithm are interleaved in parallel.

We prove the first lower bounds on cumulative memory complexity for both sequential classical computation and quantum circuits. Moreover, we develop general paradigms for bounding cumulative memory complexity inspired by the standard paradigms for proving time-space tradeoff lower bounds that can only lower bound the maximum space used during an execution. The resulting lower bounds on cumulative memory that we obtain are just as strong as the best time-space tradeoff lower bounds, which are very often known to be tight.

Although previous results for pebbling and random oracle models have yielded time-space tradeoff lower bounds larger than the cumulative memory complexity, our results show that in general computational models such separations cannot follow from known lower bound techniques and are not true for many functions.

Among many possible applications of our general methods, we show that any classical sorting algorithm with success probability at least $1/\text{poly}(n)$ requires cumulative memory $\tilde{\Omega}(n^2)$, any classical matrix multiplication algorithm requires cumulative memory $\Omega(n^6/T)$, any quantum sorting circuit requires cumulative memory $\Omega(n^3/T)$, and any quantum circuit that finds k disjoint collisions in a random function requires cumulative memory $\Omega(k^3n/T^2)$.

2012 ACM Subject Classification Theory of computation → Oracles and decision trees; Theory of computation → Quantum query complexity; Theory of computation → Quantum complexity theory

Keywords and phrases Cumulative memory complexity, time-space tradeoffs, branching programs, quantum lower bounds

Digital Object Identifier 10.4230/LIPIcs.ICALP.2023.17

Category Track A: Algorithms, Complexity and Games

Related Version *Full Version:* <https://arxiv.org/abs/2301.05680> [18]

Funding *Paul Beame:* Research supported by NSF grant CCF-2006359.

Acknowledgements Many thanks to David Soloveichik for his guidance and contributions to our initial results.

1 Introduction

For some problems, algorithms can use additional memory for faster running times or additional time to reduce memory requirements. While there are different kinds of tradeoffs between time and space, the most common complexity metric for such algorithms is the maximum time-space (TS) product. This is appropriate when a machine must allocate an



© Paul Beame and Niels Kornerup;

licensed under Creative Commons License CC-BY 4.0

50th International Colloquium on Automata, Languages, and Programming (ICALP 2023).

Editors: Kousha Etessami, Uriel Feige, and Gabriele Puppis; Article No. 17; pp. 17:1–17:20

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



algorithm’s maximum space throughout its computation. However, recent technologies like AWS Lambda [15] suggest that in the context of cloud computing, space can be allocated to a program only as it is needed. When using such services, analyzing the average memory used per step leads to a more accurate picture than measuring the maximum space.

Cumulative memory (CM), the sum over time of the space used per step of an algorithm, is an alternative notion of time-space complexity that is more fair to algorithms with rare spikes in memory. Cumulative memory complexity was introduced by Alwen and Serbinenko [12] who devised it as a way to analyze time-space tradeoffs for “memory hard functions” like password hashes. Since then, lower and upper bounds on the CM of problems in structured computational models using the black pebble game have been extensively studied, beginning with the work of [12, 7, 32, 10, 9, 8]. Structured models via pebble games are natural in the context of the random oracle assumptions that are common in cryptography. By carefully interweaving their memory-intensive steps, authors of these papers devise algorithms for cracking passwords that compute many hashes in parallel using only slightly more space than is necessary to compute a single hash. While such algorithms can use parallelism to amortize costs and circumvent proven single instance TS complexity lower bounds, their cumulative memory only scales linearly with the number of computed hashes. Strong CM results have also been shown for the black-white pebble game and used to derive related bounds for resolution proof systems [11].

The ideas used for these structured models yield provable separations between CM and TS complexity in pebbling and random oracle models. The key question that we consider is whether or not the same applies to general models of computation without cryptographic or black-box assumptions: Are existing time-space tradeoff lower bounds too pessimistic for a world where cumulative memory is more representative of a computation’s cost?

Our Results

The main answer we provide to this question is negative for both classical and quantum computation: We give *generic methods* that convert existing paradigms for obtaining time-space tradeoff lower bounds involving worst-case space to new lower bounds that replace the time-space product by cumulative space, immediately yielding a host of new lower bounds on cumulative memory complexity. With these methods, we show how to extend virtually all known proofs for time-space tradeoffs to equivalent lower bounds on cumulative memory complexity, implying that there cannot be cumulative memory savings for these problems. Our results, like those of existing time-space tradeoffs, apply in models in which arbitrary sequential computations may be performed between queries to a read-only input. Our lower bounds also apply to randomized and quantum algorithms that are allowed to make errors.

Classical computation. We focus on lower bound paradigms that apply to computations of multi-output functions $f : D^n \rightarrow R^m$. Borodin and Cook [22] introduced a method for proving time-space tradeoff lower bounds for such functions that takes a property such as the following: for some $K = K(R, n)$, constant γ , and distribution μ on D^n :

- (*) For any partial assignment τ of $k \leq \gamma m$ output values over R and any restriction (i.e., partial assignment) π of $h = h(k, n)$ coordinates on D^n ,

$$\Pr_{x \sim \mu} [f(x) \text{ is consistent with } \tau \mid x \text{ is consistent with } \pi] \leq K^{-k}.$$

and derives a lower bound of the following form:

■ **Table 1** All CM bounds match the TS lower bound when considering RAM computation or quantum circuits. The symbol * indicates that the result requires additional assumptions.

Problem	TS Lower Bound	Source	Matching CM Bound
Ranking, Sorting	$\Omega(n^2/\log n)$	[22]	Corollary 4.4
Unique Elements, Sorting	$\Omega(n^2)$	[16]	Corollary 4.14
Matrix-Vector Product (\mathbb{F})	$\Omega(n^2 \log \mathbb{F})$	[4]	Corollary 4.5
Matrix-Multiplication (\mathbb{F})	$\Omega((n^6 \log \mathbb{F})/T)$	[4]	Corollary 4.15
Hamming Closeness	$\Omega(n^{2-o(1)})$	[19]*	Full paper [18]*
Element Distinctness	$\Omega(n^{2-o(1)})$	[19]*	Full paper [18]*
Quantum Sorting	$\Omega(n^3/T)$	[29]	Theorem 3.6
Quantum k disjoint collisions	$\Omega(k^3 n/T^2)$	[27]	Corollary 4.16
Quantum Boolean Matrix-Mult	$\Omega(n^5/T)$	[29]	Full paper [18]*

► **Proposition 1.1** ([22]). *Assume that Property (*) holds for $f : D^n \rightarrow R^m$ with $\gamma > 0$ constant. Then, $T(S + \log_2 T)$ is $\Omega(m h(S/\log_2 K, n) \log K)$.*

In particular, since $S \geq \log_2 n$ is essentially always required, if we have the typical case that $h(k, n) = k^\Delta h_1(n)$ for some function $h_1(n)$ then this says that $T \cdot S^{1-\Delta}$ is $\Omega(m h_1(n) \log^{1-\Delta} K)$ or, equivalently, that $\max(S, \log n)$ is $\Omega([(m h_1(n)/T]^{1/(1-\Delta)} \log K])$. As a simplified example of our new general paradigm, we prove the following analog for cumulative complexity:

► **Theorem 1.2.** *Suppose that Property (*) holds for $f : D^n \rightarrow R^m$ with $h(k, n) = k^\Delta h_1(n)$ and $\gamma > 0$ constant. If $T \log_2 T$ is $o(m h_1(n) \log K)$ then any algorithm computing f requires cumulative memory $\Omega([(m h_1(n))^{1/(1-\Delta)} \log K] / T^{\Delta/(1-\Delta)})$.*

We note that this bound corresponds exactly to the bound on the product of time and space from Borodin-Cook method. The full version of our general theorem for randomized computation (Theorem 4.8) is inspired by an extension by Abrahamson [4] of the Borodin-Cook paradigm to average case complexity.

Our full paper ([18]) also shows how the paradigms for the best time-space tradeoff lower bounds for single-output Boolean functions, which are based on the densities of *embedded rectangles* where these functions are constant, can be extended to yield cumulative memory bounds.

Quantum computation. We develop an extension of our general approach that applies to quantum computation as well. In this case Property (*) and its extensions that we use for our more general theorem must be replaced by statements about quantum circuits with a small number of queries. In this case, we first generalize the quantum time-space tradeoff for sorting proven in [29], which requires that the time order in which output values are produced must correspond to the sorted order, to a matching cumulative memory complexity bound of $\Omega(n^3/T)$ that works for any fixed time-ordering of output production, yielding a more general lower bound. (For example, an algorithm may be able to determine the median output long before it determines the other outputs.) We then show how an analog of our classical general theorem can be applied to extend to paradigms for quantum time-space tradeoffs to cumulative memory complexity bounds for other problems.

A summary of our results for both classical and quantum complexity is given in Table 1.

Previous work

Memory hard functions and cumulative memory complexity. Alwen and Serbinenko [12] introduced parallel cumulative (memory) complexity as a metric for analyzing the space footprint required to compute *memory hard functions (MHFs)*, which are functions designed to require large space to compute. Most MHFs are constructed using hashgraphs [26] of DAGs whose output is a fixed length string and their proofs of security are based on pebbling arguments on these DAGs while assuming access to truly random hash functions for their complexity bounds [12, 21, 32, 8, 10, 20]. (Also see our full paper [18] for their use in separating CM and TS complexity.) Recent constructions do not require random hash functions; however, they still rely on cryptographic assumptions [25, 14].

Classical time-space tradeoffs. While these were originally studied in restricted pebbling models similar to those considered to date for cumulative memory complexity [35, 23], the gold-standard model for time-space tradeoff analysis is that of unrestricted branching programs, which simultaneously capture time and space for general sequential computation. Following the methodology of Borodin and Cook [22], who proved lower bounds for sorting, many other problems have been analyzed (e.g., [37, 2, 3, 16, 30]), including universal hashing and many problems in linear algebra [4]. (See [34, Chapter 10] for an overview.) A separate methodology for single-output functions, introduced in the context of restricted branching programs [24, 31], was extended to general branching programs in [17], with further applications to other problems [5] including multi-precision integer multiplication [33] and error-correcting codes [28] as well as over Boolean input domains [6, 19]. Both of these methods involve breaking the program into blocks to analyze the computation under natural distributions over the inputs based on what happens at the boundaries between blocks.

Quantum time-space tradeoffs. Similar blocking strategies can be applied to quantum circuits to achieve time-space trade-offs for multi-output functions. In [29] the authors use direct product theorems to prove time-space tradeoffs for sorting and Boolean matrix multiplication. They also proved somewhat weaker lower bounds for computing matrix-vector products for fixed matrices A ; those bounds were extended in [13] to systems of linear inequalities. However, both of these latter results apply to computations where the fixed matrix A defining the problem depends on the space bound and, unlike the case of sorting or Boolean matrix multiplication, do not yield a fixed problem for which the lower bound applies at all space bounds. More recently [27] extended the recording query technique of Zhandry in [38] to obtain time-space lower bounds for the k -collision problem and match the aforementioned result for sorting.

Our methods

At the highest level, we employ part of the same paradigms previously used for time-space tradeoff lower bounds. Namely breaking up the computations into blocks of time and analyzing properties of the branching programs or quantum circuits based on what happens at the boundaries between time blocks. However, for cumulative memory complexity, those boundaries cannot be at fixed locations in time and their selection needs to depend on the space used in these time steps.

Further, in many cases, the time-space tradeoff lower bound needs to set the lengths of those time blocks in a way that depends on the specific space bound. When extending the ideas to bound cumulative memory usage, there is no single space bound that can be used

throughout the computation; this sets up a tricky interplay between the choices of boundaries between time blocks and the lengths of the time blocks. Because the space usage within a block may grow and shrink radically, even with optimal selection of block boundaries, the contribution of each time block to the overall cumulative memory may be significantly lower than the time-space product lower bound one would obtain for the individual block.

We show how to bound any loss in going from time-space tradeoff lower bounds to cumulative memory lower bounds in a way that depends solely on the bound on the lengths of blocks as a function h_0 of the target space bound (cf. Lemma 4.7). For many classes of bounding functions we are able to bound the loss by a constant factor, and we are able to show that it is always at most an $O(\log n)$ factor loss. If this bounding function h_0 is non-constant, we also need to bound the optimum way for the algorithm to allocate its space budget for producing the required outputs throughout its computation. This optimization again depends on the bounding function h_0 . This involves minimizing a convex function based on h_0 subject to a mix of convex and concave constraints, which is not generally tractable. However, assuming that h_0 is nicely behaved, we are able to apply specialized convexity arguments (cf. Lemma 4.10) which let us derive strong lower bounds on cumulative memory complexity.

Road map. We give the overall definitions in Section 2, including a review of the standard definitions of the work space used by quantum circuits. In Section 3, we give our lower bound for quantum sorting algorithms which gives a taste of the issues involved for our general theorems. In Section 4, we give the general theorems that let us convert the Borodin-Cook-Abrahamson paradigm for multi-output functions to cumulative memory lower bounds for classical randomized algorithms; that section also contains the corresponding theorems for quantum lower bounds and statements of some sample applications for our general results. Appendix A contains the arguments that bound the optimum allocations of cumulative space budgets to time steps. Our full paper [18] contains more details, a conditional separation between CM and TS complexity, detailed applications of the general theorems we present here, and our bounds for single-output functions.

2 Preliminaries

Cumulative memory is an abstract notion of time-space complexity that can be applied to any model of computation with a natural notion of space. Here we will use branching programs and quantum circuits as concrete models, although our results generalize to any reasonable model of computation.

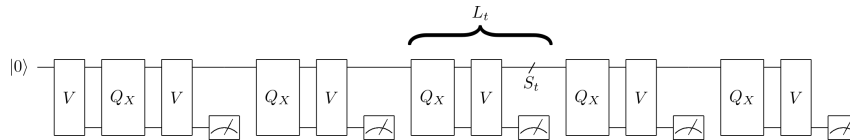
Branching Programs. A branching program with input $\{x_1, \dots, x_n\} \in D^n$ is defined using a rooted DAG in which each non-sink vertex is labeled with an $i \in [n]$ and has $|D|$ outgoing edges that correspond to possible values of x_i . Each edge is optionally labeled by some number of output statements expressed as pairs (j, o_j) where $j \in [m]$ is an output index and $o_j \in R$ (if outputs are to be ordered) or simply $o_j \in R$ (if outputs are to be unordered). Evaluation starts at the root v_0 and follows the appropriate labels of the respective x_i . We consider branching programs P that contain $T + 1$ layers where the outgoing edges from nodes in each layer t are all in layer $t + 1$. We impose no restriction on the query pattern of the branching program or when it can produce parts of the output. The *time* of the branching program is $T(P) = T$. The *space* of the branching program is $S(P) = \max_t \log_2 |L_t|$ where L_t is the set of nodes in layer t . Observe that in the absence of any limit on its space, a branching program could be a decision tree; hence the minimum time for branching programs

17:6 Computational Cumulative Memory Lower Bounds

to compute a function f is its *decision tree complexity*. The *time-space* (product) used by the branching program is $TS(P) = T(P)S(P)$. The *cumulative memory* used by the branching program is $CM(P) = \sum_t \log_2 |L_t|$.

Branching programs are very general, and simultaneously model time and space for sequential computation. In particular they model time and space for random-access off-line multitape Turing machines and random-access machines (RAMs) when time is unit-cost, space is log-cost, and the input and output are read-only and write-only respectively. Branching programs are much more flexible than these models since they can make arbitrary changes to their storage in a single step.

Quantum Circuits. We also consider quantum circuits \mathcal{C} classical read-only input $X = x_1, \dots, x_n$ that can be queried using an XOR query oracle. As is normal in circuit models, each output wire is associated with a fixed position in the output sequence, independent of the input. As shown in Figure 1 following [29], we abstract an arbitrary quantum circuit \mathcal{C} into layers $\mathcal{C} = \{L_1, \dots, L_T\}$ where layer L_t starts with the t -th query Q to the input and ends with the start of the next layer. During each layer, an arbitrary unitary transformation V gets applied which can express an arbitrary sub-circuit involving input-independent computation. The sub-circuit/transformation V outputs S_t qubits for use in the next layer in addition to some qubits that are immediately measured in the standard basis, some of which are treated as classical write-only output. The time of \mathcal{C} is lower bounded by the number of layers T and we say that the space of layer L_t is S_t . Observe that to compute a function f , T must be at least the *quantum query complexity* of f since that measure corresponds the above circuit model when the space is unbounded. Note that the cumulative memory of a circuit is lower-bounded by the sum of the S_t . For convenience we define S_0 , the space of the circuit before its first query, to be zero. Thus we only consider the space after the input is queried.



■ **Figure 1** The abstraction of a quantum circuit into layers.

3 Quantum cumulative memory complexity of sorting

As an illustrative example, we first show that the quantum cumulative memory complexity of sorting is $\Omega(n^3/T)$, matching the TS complexity bounds given in [29, 27]. This involves the quantum circuit model which, as we have noted, produces each output position at a predetermined input-independent layer. We restrict our attention to circuits that output all elements in the input in some fixed rank order. While our proof is inspired by the time-space lower bound of [29], it can be easily adapted to follow the proof in [27] instead. We start by constructing a probabilistic reduction from the k -threshold problem to sorting.

► **Definition 3.1.** *In the k -threshold problem we receive an input $X = x_1, \dots, x_n$ where $x_i \in \{0, 1\}$. We want to accept iff there are at least k distinct values for i where $x_i = 1$.*

► **Proposition 3.2** (Theorem 13 in [29]). *For every $\gamma > 0$ there is an $\alpha > 0$ such that any quantum k -threshold circuit with at most $T \leq \alpha\sqrt{kn}$ queries and with perfect soundness must have completeness $\sigma \leq e^{-\gamma k}$ on inputs with Hamming weight k .*

► **Lemma 3.3.** *Let $\gamma > 0$. Let n be sufficiently large and $\mathcal{C}(X)$ be a quantum circuit with input $X = x_1, \dots, x_n$. There is a $\beta < 1$ depending only on γ such that for all $k \leq \beta^2 n$ and $R \subseteq \{n/2 + 1, \dots, n\}$ where $|R| = k$, if $\mathcal{C}(X)$ makes at most $\beta\sqrt{kn}$ queries, then the probability that $\mathcal{C}(X)$ can correctly output all k pairs (x_i, r_j) where $r_j \in R$ and x_i is the r_j -th smallest element of X is at most $e^{(1-\gamma)k-1}$. If R is a contiguous set of integers, then the probability is at most $e^{-\gamma k}$.*

A version of this lemma was first proved in [29] with the additional assumption that the set of output ranks R is a contiguous set of integers; this was sufficient to show that any quantum circuit that produces its sorted output in sorted time order requires that T^2S is $\Omega(n^3)$. The authors stated that their proof can be generalized to any fixed rank ordering, but the generalization is not obvious. We generalize their lemma to non-contiguous R , which is sufficient to obtain an $\Omega(n^3/T)$ lower bound on the cumulative complexity of sorting independent of the time order in which the sorted output is produced.

Proof of Lemma 3.3. Choose α as the constant for γ in Proposition 3.2 and let $\beta = \sqrt{2}\alpha/6$. Let \mathcal{C} be a circuit with at most $\beta\sqrt{kn}$ layers that outputs the k correct pairs (x_i, r_j) with probability p . Let $R = \{r_1, \dots, r_k\}$ where $r_1 < r_2 < \dots < r_k$. We describe our construction of a circuit $\mathcal{C}'(X)$ solving the k -threshold problem on inputs $X = x_1, \dots, x_{n/2}$ with exactly k ones in terms of a function $f : [n/2] \rightarrow R$. Given f , we re-interpret the input as follows: we replace each x_i with $x'_i = f(i)x_i$, add k dummy values of 0, and add one dummy value of j for each $j \in \{n/2 + 1, \dots, n\} \setminus R$. Doing this gives us an input $X' = x'_1, \dots, x'_n$ that has $n/2$ zeroes. If we assume that f is 1-1 on the k ones of X , then the image of the ones of X will be R and there will be precisely one element of X' for each $j \in \{n/2 + 1, \dots, n\}$. Therefore the element of rank $j > n/2$ in X' will have value j , and hence the rank r_1, \dots, r_k elements of X' will be the images of precisely those elements of X with $x_i = 1$.

To obtain perfect soundness, we cannot rely on the output of $\mathcal{C}(X')$ and must be able to check that each of the output ranks was truly mapped to by a distinct one of X . For each element x_i of X we simply append its index i as $\log_2 n$ low order bits to its image x'_i and append an all-zero bit-vector of length $\log_2 n$ to each dummy value to obtain input X'' . Doing so will not change the ranks of the elements in X' , but will allow recovery of the k indices that should be the ones in X . In particular, circuit $\mathcal{C}'(X)$ will run $\mathcal{C}(X'')$ and then for each output x''_j with low order bits i , $\mathcal{C}'(X)$ will query x_i , accepting if and only if all of those $x_i = 1$. More precisely, since the mapping from each x_i to the corresponding x''_i is only a function of f , x_i , and i , as long as $\mathcal{C}'(X)$ has an explicit representation of f , it can simulate each query of $\mathcal{C}(X'')$ with two oracle queries to X . Since \mathcal{C}' has at most

$$2\beta\sqrt{kn} + k \leq 3\beta\sqrt{kn} \leq \alpha\sqrt{kn/2}$$

layers, by Proposition 3.2, it can only accept with probability $\leq e^{-\gamma k}$ on inputs with k ones.

We now observe that for each fixed X with exactly k ones, for a randomly chosen function $f : [n/2] \rightarrow R$, the probability that f is 1-1 on the ones of X' is exactly $k!/k^k \geq e^{1-k}$. Therefore $\mathcal{C}'(X)$ will give the indices of the k ones in X with probability¹ at least $p \cdot e^{1-k}$. However, this probability must be at most $e^{-\gamma k}$, so we can conclude that $p \leq e^{(1-\gamma)k-1}$. In the event that R is a contiguous set of integers, observe that any choice for the function f will make X'' have the ones of X become ranks r_1, \dots, r_k . So the probability of finding the ones is at least $p \leq e^{-\gamma k}$. ◀

¹ Note that though this is exponentially small in k it is still sufficiently large compared to the completeness required in the lower bound for the k -threshold problem.

17:8 Computational Cumulative Memory Lower Bounds

By setting k and γ appropriately, Lemma 3.3 gives a useful upper bound on the number of fixed ranks successfully output by any $\beta\sqrt{Sn}$ query quantum circuit that has access to S qubits of input dependent initial state. To handle input-dependent initial state, we will need to use the following proposition.

► **Proposition 3.4** ([1]). *Let \mathcal{C} be a quantum circuit, ρ be any S qubit (possibly mixed) state, and I be the S qubit maximally mixed state. If \mathcal{C} with initial state ρ produces some output \mathcal{O} with probability p , then \mathcal{C} with initial state I produces \mathcal{O} with probability at least $p/2^{2S}$.*

This allows us to bound the overall progress made by any short quantum circuit.

► **Lemma 3.5.** *There is a constant $\beta > 0$ such that, for any fixed set of $S \leq \beta^2 n$ ranks that are greater than $n/2$, the probability that any quantum circuit \mathcal{C} with at most $\beta\sqrt{Sn}$ queries and S qubits of input-dependent initial state correctly produces the outputs for these S ranks is at most $1/e$.*

Proof. Choose β as the constant when γ is $1 + \ln(4)$ in Lemma 3.3. Applying Proposition 3.4 to the bound in Lemma 3.3 gives us that a quantum circuit with S qubits of input-dependent state can produce a fixed set of $k \leq \beta^2 n$ outputs larger than median with a probability at most $2^{2S} e^{(1-\gamma)k-1}$. Since $\gamma = 1 + \ln(4)$ setting $k = S$ gives that this probability is $\leq 1/e$. ◀

► **Theorem 3.6.** *When n is sufficiently large, any quantum circuit \mathcal{C} for sorting a list of length n with success probability at least $1/e$ and at most T layers that produces its sorted outputs in any fixed time order requires cumulative memory that is $\Omega(n^3/T)$.*

Proof. We partition \mathcal{C} into blocks with large cumulative memory that can only produce a small number of outputs. We achieve this by starting at last unpartitioned layer and finding a suitably low space layer before it so that we can apply Lemma 3.5 to upper bound the number of correct outputs that can be produced in that block with a success probability of at least $1/e$. Let β be the constant from Lemma 3.5 and $k^*(t)$ be the least non-negative integer value of k such that the interval:

$$I(k, t) = \left[t - \frac{\beta}{2}(2^{k+1} - 1)\sqrt{n}, t - \frac{\beta}{2}(2^k - 1)\sqrt{n} \right]$$

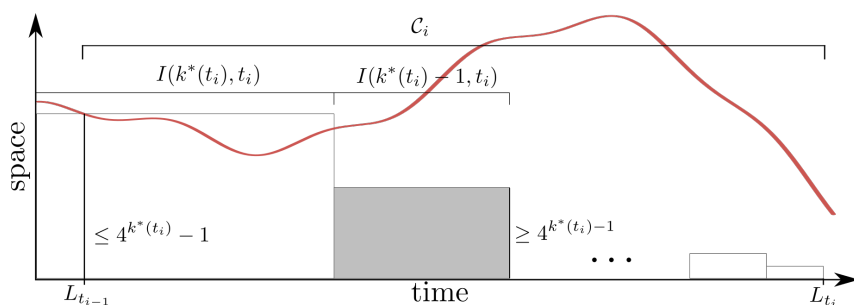
contains some t' such that $S_{t'} \leq 4^k - 1$. We recursively define our blocks as follows. Let ℓ be the number of blocks generated by this method. The final block \mathcal{C}_ℓ starts with the first layer $t_{\ell-1} \in I(k^*(T), T)$ where $S_{t_{\ell-1}} \leq 4^{k^*(T)} - 1$ and ends with layer $t_\ell = T$. Let t_i be the first layer of block \mathcal{C}_{i+1} . Then the block \mathcal{C}_i starts with the first layer $t_{i-1} \in I(k^*(t_i), t_i)$ where $S_{t_{i-1}} \leq 4^{k^*(t_i)} - 1$ and ends with t_i . See Figure 2 for an illustration of our partitioning. Since $S_0 = 0$ we know that $k^*(t) \leq \log(T)$. Likewise since $S_t > 0$ when $t > 0$, for all $t > \frac{\beta}{2}\sqrt{n}$ we know that $0 < k^*(t) \leq \log(T)$.

Block \mathcal{C}_i starts with less than $4^{k^*(t_i)}$ qubits of initial state and has length at most $\beta 2^{k^*(t_i)}\sqrt{n}$; so by Lemma 3.5, if $4^{k^*(t_i)} \leq \beta^2 n$, the block \mathcal{C}_i can output at most $4^{k^*(t_i)}$ inputs with failure probability at most $1/e$. Additionally \mathcal{C}_i has at least $\frac{\beta}{2} 2^{k^*(t_i)-1}\sqrt{n}$ layers so

$$\sum_{i=1}^{\ell} \frac{\beta}{4} 2^{k^*(t_i)}\sqrt{n} \leq T \tag{1}$$

and each of these layers has at least $4^{k^*(t_i)-1}$ qubits², so the cumulative memory of \mathcal{C}_i is at

² This may not hold for \mathcal{C}_1 with length less than $\frac{\beta}{2}\sqrt{n}$, but Lemma 3.3 gives us that this number of layers is insufficient to find a fixed rank input with probability at least $1/e$. Thus we can omit such a block from our analysis.



■ **Figure 2** How we define the block C_i that ends at layer L_{t_i} . The red line is a plot of C_i 's space over time. The grey layers are the ones used to lower bound the cumulative memory complexity of C_i , as each of these layers uses at least $4^{k^*(t_i)-1}$ qubits and the length of this interval is $\frac{\beta}{2}2^{k^*(t_i)-1}\sqrt{n}$.

least $\frac{\beta}{2}2^{3k^*(t_i)-3}\sqrt{n}$ so

$$CM(C) \geq \sum_{i=1}^{\ell} \frac{\beta}{2} 2^{3k^*(t_i)-3} \sqrt{n}. \quad (2)$$

We now have two possibilities: If we have some i such that $4^{k^*(t_i)} > \beta^2 n$, the cumulative memory of C_i alone is at least $\beta^4 n^2 / 16$ which is $\Omega(n^2)$ and hence C has cumulatively memory $\Omega(n^3/T)$ since $T \geq n$. Otherwise, since we require that the algorithm is correct with probability at least $1/e$, each block C_i can produce at most $4^{k^*(t_i)}$ outputs. Since our circuit must output all $n/2$ elements larger than the median, we know $\sum_{i=1}^{\ell} 4^{k^*(t_i)} \geq n/2$. For convenience we define $w_i = 2^{k^*(t_i)}$ which allows us to express the constraints as

$$CM(C) \geq \frac{\beta}{16} \sqrt{n} \sum_{i=1}^{\ell} w_i^3 \quad \text{and} \quad \frac{\beta}{4} \sqrt{n} \sum_{i=1}^{\ell} w_i \leq T \quad \text{and} \quad \sum_{i=1}^{\ell} w_i^2 \geq n/2. \quad (3)$$

Minimizing $\sum_{i=1}^{\ell} w_i^3$ is a non-convex optimization problem and can instead be solved using

$$\text{Minimize} \quad \sum_{i=1}^{\ell} x_i^3 \quad \text{subject to} \quad \sum_{i=1}^{\ell} x_i^2 \geq \xi \quad \text{and} \quad \sum_{i=1}^{\ell} x_i \leq \xi \quad \text{and} \quad \forall i, x_i \geq 0, \quad (4)$$

for $x_i = \frac{8T}{\beta n^{3/2}} w_i$ and $\xi = \frac{32T^2}{\beta^2 n^2}$. Lemma A.1 from Appendix C shows that for non-negative x_i with $\sum x_i \leq \sum x_i^2$, we have $\sum x_i^2 \leq \sum x_i^3$. Thus $\sum x_i^3 \geq \xi$ and applying the variable substitution gives us: $\sum_{i=1}^{\ell} w_i^3 \geq \frac{\beta n^{5/2}}{16T}$. Plugging this into Equation (3) gives us the bound:

$$CM(C) \geq \frac{\beta^2 n^3}{256T} \quad \text{and hence the cumulative memory of } C \text{ is } \Omega(n^3/T). \quad \blacktriangleleft$$

4 General methods for proving cumulative memory lower bounds

Our method involves adapting techniques previously used to prove tradeoff lower bounds on worst-case time and worst-case space. We show that the same properties that yield lower bounds on the product of time and space in the worst case can also be used to produce nearly identical lower bounds on cumulative memory. To do so, we first revisit the standard approach to such time-space tradeoff lower bounds.

The standard method for time-space tradeoff lower bounds for multi-output functions

Consider a multi-output function f on D^n where the output $f(x)$ is either unordered (the output is simply a set of elements from R) or ordered (the output is a vector of elements from R). Then $|f(x)|$ is either the size of the set or the length of the vector of elements. The standard method for obtaining an ordinary time-space tradeoff lower bounds for multi-output functions on D -way branching programs is the following:

The part that depends on f . Choose a suitable probability distribution μ on D^n , often simply the uniform distribution on D^n and then:

(A) Prove that $\Pr_{x \sim \mu}[|f(x)| \geq m] \geq \alpha$.

(B) Prove that for all $k \leq m'$ and any branching program B of height $\leq h'(k, n)$, the probability for $x \sim \mu$ that B produces at least k correct output values of f on input x is at most $C \cdot K^{-k}$ for some $m', h', K = K(R, n)$, and constant C independent of n .

Observe that under any distribution μ , a branching program with ordered outputs that makes no queries can produce k outputs that are all correct with probability at least $|R|^{-k}$, so the bound in (B) shows that, roughly, up to the difference between K and $|R|$ there is not much gained by using a branching program of height h .

The generic completion. In the following outline we omit integer rounding for readability.

■ Let $S' = S + \log_2 T$ and suppose that

$$S' \leq m' \log_2 K - \log_2(2C/\alpha). \quad (5)$$

■ Let $k = \lfloor S' + \log_2(2C/\alpha) \rfloor / \log_2 K$, which is at most m' by hypothesis on S' , and define $h(S', n) = h'(k, n)$.

■ Divide time T into $\ell = T/h$ blocks of length $h = h(S', n)$.

■ The original branching program can be split into at most $T \cdot 2^S = 2^{S'}$ sub-branching programs of height $\leq h$, each beginning at a boundary node between layers. By Property (B) and a union bound, for $x \sim \mu$ the probability that at least one of these $\leq 2^{S'}$ sub-branching programs of height at most h produces k correct outputs on input x is at most $2^{S'} \cdot C \cdot K^{-k} \leq \alpha/2$ by our choice of k .

■ Under distribution μ , by (A), with probability at least α , an input $x \sim \mu$ has some block of time where at least $m/\ell = m \cdot h(S', n)/T$ outputs of f must be produced on input x .

■ If $m \cdot h(S', n)/T \leq k$, this can occur for at most an $\alpha/2$ fraction of inputs under μ . Therefore we have $m \cdot h(S', n)/T > k = \lfloor S' + \log_2(2C/\alpha) \rfloor / \log_2 K$ and hence since $h(S', n) \geq h(S, n)$, combining with Equation (5), we have

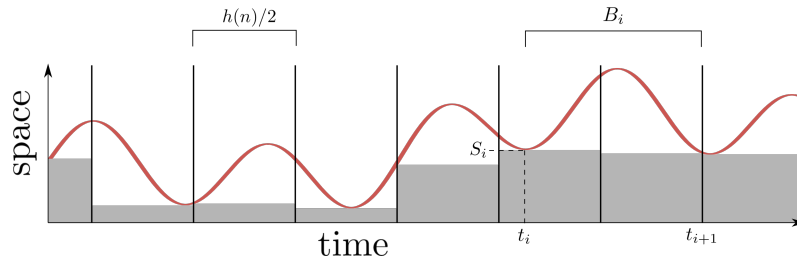
$$T \cdot (S + \log_2 T) = T \cdot S' \geq \min(m h(S, n), m' n') \log_2 K - \log_2(C/\alpha) \cdot T$$

where $n' \leq n$ is the decision tree complexity of f and hence a lower bound on T .

► **Remark 4.1.** Though it will not impact our argument, for many instances of the above outline, the proof of Property (B) is shown for a decision tree of the same height by proving an analog for the conditional probability along each path in the decision tree separately; this will apply to the tree as a whole since the paths are followed by disjoint inputs, so Property (B) follows from the alternative property below:

(B') For any partial assignment τ of $k \leq m'$ output values over R and any restriction (i.e., partial assignment) π of $h'(k, n)$ coordinates within D^n ,

$$\Pr_{x \sim \mu} [f(x) \text{ is consistent with } \tau \mid x \text{ is consistent with } \pi] \leq C \cdot K^{-k}.$$



■ **Figure 3** Our generic method for choosing blocks when $h(k, n) = h(n)$. The area marked in grey corresponds to the cumulative memory lower bound we obtain.

Observe that Property (B') is only a slightly more general version of Property (*) from the introduction where $C = 1$, m' is arbitrary, and h' is used instead of h .

► **Remark 4.2.** The above method still gives lower bounds for many multi-output functions $g : D^N \rightarrow R^M$ that have individual output values that are easy to compute or large portions of the input space on which they are easy to compute. The bounds follow by applying the method to some subfunction f of g given by $f(x) = \Pi_O(g(x, \pi))$ where π is a partial assignment to the input coordinates and Π_O is a projection onto a subset O of output coordinates. In the subsequent discussions we ignore this issue, but the idea can be applied to all of our lower bound methods.

A general extension to cumulative memory bounds

To give a feel for the basic ideas of the method, we first show this for a simple case. Observe that, other than the separate bound on time, the lower bound on cumulative memory usage we prove in this case is asymptotically identical to the bound achieved for the product of time and worst-case space using the standard outline.

► **Theorem 4.3.** *Let $c > 0$. Suppose that properties (A) and (B) apply for $h'(k, n) = h(n)$, $m' = m$, and $\alpha = C = 1$. If $T \log_2 T \leq \frac{m h(n) \log_2 K}{6(c+1)}$ then the cumulative memory used in computing $f : D^n \rightarrow R^m$ in time T with success probability at least T^{-c} is at least $\frac{1}{6} m h(n) \log_2 K$.*

Proof. Fix a deterministic branching program P of length T computing f . Rather than choosing fixed blocks of height $h = h(n)$, layers of nodes at a fixed distance from each other, and a fixed target of k outputs per block, we choose the block boundaries depending on the properties of P and the target k depending on the property of the boundary layer chosen.

Let $H = \lfloor h(n)/2 \rfloor$. We break P into $\ell = \lceil T/H \rceil$ time segments of length H working backwards from step T so that the first segment may be shorter than the rest. We let $t_1 = 0$ and for $1 < i \leq \ell$ we let $t_i = \arg \min \{ |L_t| : T - (\ell - i + 1) \cdot H \leq t < T - (\ell - i) \cdot H \}$ be the time step with the fewest nodes among all time steps $t \in [T - (\ell - i + 1) \cdot H, T - (\ell - i) \cdot H]$.

The i -th time block of P will be between times t_i and t_{i+1} . Observe that by construction $|t_{i+1} - t_i| \leq h(n)$ so each block has length at most $h(n)$. This construction is shown in Figure 3 Set $S_i = \log_2 |L_{t_i}|$ so that L_{t_i} has at 2^{S_i} nodes. By definition of each t_i , the cumulative memory used by P ,

$$CM(P) \geq \sum_{i=1}^{\ell} S_i \cdot H. \tag{6}$$

(Note that since $S_1 = 0$, it does not matter that the first segment is shorter than the rest³.)

³ This simplifies some calculations and is the prime reason for starting the time segment boundaries at T .

17:12 Computational Cumulative Memory Lower Bounds

We now define the target k_i for the number of output values produced in each time block to be the smallest integer such that $K^{-k_i} \leq 2^{-S_i}/T^{c+1}$. That is,

$$k_i = \lceil (S_i + (c+1) \log_2 T) / \log_2 K \rceil.$$

For $x \sim \mu$, for each $i \in [\ell]$ and each sub-branching program B rooted at some node in L_{t_i} and extending until time t_{i+1} , by our choice of k_i and Property (B), if $k_i \leq m$, the probability that B produces at least k_i correct outputs on input x is at most $2^{-S_i}/T^{c+1}$. Therefore, by a union bound, for $x \sim \mu$ the probability that P produces at least k_i correct outputs in the i -th time block on input x is at most $|L_{t_i}| \cdot 2^{-S_i}/T^{c+1} = 1/T^{c+1}$. Therefore, if each $k_i \leq m$, the probability for $x \sim \mu$ that there is some i such that P produces at least k_i correct outputs on input x during the i -th block is at most $\ell/T^{c+1} < T^c$ and the probability for $x \sim \mu$ that P produces at most $\sum_{i=1}^{\ell} (k_i - 1)$ correct outputs in total on input x is $> 1 - 1/T^c$.

If each $k_i \leq m$, since P must produce m correct outputs on $x \in D^n$ with probability at least $1/T^c$, we must have $\sum_{i=1}^{\ell} (k_i - 1) \geq m$. On the other hand, if some $k_i > m$ we have the same bound. Using our definition of k_i we have $\sum_{i=1}^{\ell} \lceil (S_i + (c+1) \log_2 T) / \log_2 K \rceil \geq m$ or $\sum_{i=1}^{\ell} (S_i + (c+1) \log_2 T) \geq m \cdot \log_2 K$. Plugging in the bound (6) on the cumulative memory and the value of ℓ , it implies that $CM(P)/H + (c+1) \lceil T/H \rceil \cdot \log_2 T \geq m \cdot \log_2 K$ or that $CM(P) + (c+1)T \log_2 T \geq \frac{1}{3} m \cdot h(n) \cdot \log_2 K$, where the 3 on the right rather than a 2 allows us to remove the ceiling. Therefore either

$$T \log_2 T > \frac{m \cdot h(n) \cdot \log_2 K}{6(c+1)} \quad \text{or} \quad CM(P) \geq \frac{1}{6} m h(n) \log_2 K. \quad \blacktriangleleft$$

Simple applications. This simple case of our general theorem is sufficient to obtain many tight new lower bounds on cumulative memory complexity including the following (full proofs are in [18]):

► **Corollary 4.4.** *Producing the ranks (positions of each input in the sorted order) or sorting n integers from $[n^2]$ requires CMC that is $\Omega(n^2 / \log_2 n)$.*

► **Corollary 4.5.** *For many fixed (random or explicit) $n \times n$ matrices A , computing Ax over a finite field \mathbb{F} requires CMC that is $\Omega(n^2 \log |\mathbb{F}|)$.*

Corollary 4.4 uses property (B) for ranking with $m' = h'(k, n) = \Theta(n)$, $C = 1$, $K = 2^{\Theta(1/\log n)}$ proven in [22, Lemma 1]. Corollary 4.5 uses property (B') with $m' = h'(k, n) = cn$, for $0 < c \leq 1/2$, $C = 1$, and $K = |\mathbb{F}|^c$ proven in [4, Theorem 4.6].

Full general theorem. In the general version of our theorem there are a number of additional complications, most especially because the branching program height limit $h(k, n)$ in Property (B) can depend on k , the target for the number of outputs produced. This forces the lengths of the blocks and the space used at the boundaries between blocks to depend on each other in a quite delicate way. In order to discuss the impact of that dependence and state our general theorem, we need the following definition.

► **Definition 4.6.** *Given a non-decreasing function $p : \mathbb{R} \rightarrow \mathbb{R}$ with $p(1) = 1$, we define $p^{-1} : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ by $p^{-1}(R) = \min\{j \mid p(j) \geq R\}$. We also define the loss, \mathcal{L}_p , of p by*

$$\mathcal{L}_p(n) = \min_{1 \leq k \leq p(n)} \frac{\sum_{j=1}^k p^{-1}(j)}{k \cdot p^{-1}(k)}.$$

► **Lemma 4.7.** *The following hold for every non-decreasing function $p : \mathbb{R} \rightarrow \mathbb{R}$ with $p(1) = 1$:*

- (a) $1/p(n) \leq \mathcal{L}_p(n) \leq 1$.
- (b) *If p is a polynomial function $p(s) = s^{1/c}$ then $\mathcal{L}_p(n) > 1/2^{c+1}$.*
- (c) *For any $c > 1$, $\mathcal{L}_p(n) \geq \min_{1 \leq s \leq n} \frac{p(s) - p(s/c)}{cp(s)}$.*
- (d) *We say that p is nice if it is differentiable and there is an integer $c > 1$ such that for all x , $p'(cx) \geq p'(x)/c$. If p is nice then $\mathcal{L}_p(n)$ is $\Omega(1/\log_2 n)$. This is tight for p with $p(s) = 1 + \log_2 s$.*

We prove this technical lemma in the full paper [18]. Here is our full general theorem.

► **Theorem 4.8.** *Let $c > 0$. Suppose that function f defined on D^n has properties (A) and (B) with α that is $1/n^{O(1)}$ and m' that is $\omega(\log_2 n)$. For $s > 0$, define $h(s, n)$ to be $h'(k, n)$ for $k = s/\log_2 K$. Suppose that $h(s, n) = h_0(s) h_1(n)$ with $h_0(1) = 1$ and h_0 is constant or a differentiable function such that $s/h_0(s)$ is increasing and concave. Define $S^* = S^*(T, n)$ by*

$$S^*/h_0(S^*) = (m h_1(n) \log_2 K) / 6T.$$

- (a) *Either $T \log_2(2CT^{c+1}/\alpha) > \frac{1}{6} m h_1(n) \log_2 K$, which implies that T is $\Omega(\frac{m h_1(n) \log K}{\log n})$, or the cumulative memory used by a randomized branching program in computing f in time T with error $\varepsilon \leq \alpha(1 - 1/(2T^c))$ is at least*

$$\frac{1}{6} \mathcal{L}_{h_0}(n \log_2 |D|) \cdot \min(m h(S^*(T, n), n), 3m' h'(m'/2, n)) \cdot \log_2 K.$$

- (b) *Further any randomized random-access machine computing f in time T with error $\varepsilon \leq \alpha(1 - 1/(2T^c))$ requires cumulative memory*

$$\Omega(\mathcal{L}_{h_0}(n \log_2 |D|) \cdot \min(m h(S^*(T, n), n), m' h'(m'/2, n)) \cdot \log_2 K).$$

Before we give the proof of the theorem, we note that by Lemma 4.7, in the case that h_0 is constant or $h_0(s) = s^\Delta$ for some constant $\Delta > 0$, which together account for all existing applications we are aware of, the function \mathcal{L}_{h_0} is lower bounded by a constant. In the latter case, h_0 is differentiable, has $h_0(s) = 1$, and the function $s/h_0(s) = s^{1-\Delta}$ is increasing and concave so it satisfies the conditions of our theorem. By using $\alpha = 1$, $m' = m$, and $C = 1$ with h from Property (*) in place of h' in Property (B'), Theorem 4.8 yields Theorem 1.2.

More generally, the value S^* in the statement of this theorem is at least a constant factor times the value of S used in the generic time-space tradeoff lower bound methodology. Therefore, the cumulative memory lower bound in Theorem 4.8 for random-access machines is close to the lower bound on the product of time and space using standard methods.

Proof of Theorem 4.8. We prove both (a) and (b) directly for branching programs, which can model random-access machines, and will describe the small variation that occurs in the case that the branching program in question comes from a random-access machine. To prove these properties for randomized branching programs, by Yao's Lemma [36] it suffices to prove the properties for deterministic branching programs that have error at most ε under distribution μ . Fix a (deterministic) branching program P of length T computing f with error at most ε under distribution μ . Without loss of generality, P has maximum space usage at most $S^{max} = n \log_2 |D|$ space since there are at most $|D^n|$ inputs.

Let $H = \lfloor h_1(n)/2 \rfloor$. We break P into $\ell = \lceil T/H \rceil$ time segments of length H working backwards from step T so that the first segment may be shorter than the rest. We then choose a sequence of *candidates* for the time steps in which to begin new blocks, as follows: We let $\tau_1 = 0$ and for $1 < i \leq \ell$ we let

$$\tau_i = \arg \min \{ |L_t| : T - (\ell - i + 1) \cdot H \leq t < T - (\ell - i) \cdot H \}$$

17:14 Computational Cumulative Memory Lower Bounds

be the time step with the fewest nodes among all time steps $t \in [T - (\ell - i + 1) \cdot H, T - (\ell - i) \cdot H]$. Set $\sigma_i = \log_2 |L_{\tau_i}|$ so that L_{τ_i} has at 2^{σ_i} nodes. This segment contributes at least $\sigma_i \cdot H$ to the cumulative memory bound of P .

To choose the beginning t_{i^*} of the last time block⁴, we find the smallest k such that $h_0(\sigma_{\ell-k+1}) < k$. Such a k must exist since h_0 is a non-decreasing non-negative function, $h_0(1) = 1$ and $\sigma_1 = 0 < 1$. We now observe that the length of the last block is at most $k \cdot H$ which by choice of k is less than $h(\sigma_{\ell-k+1}, n)$ and hence we have satisfied the requirements for Property (B) to apply at each starting node of the last time block.

By our choice of each τ_i , the cumulative memory used in the last k segments is at least $\sum_{j=1}^k \sigma_{\ell+1-j} \cdot H$. Further, since k was chosen as smallest with the above property, we know that for every $j \in [k-1]$ we have $h_0(\sigma_{\ell-j+1}) \geq j$. Hence we have $\sigma_{\ell-j+1} \geq h_0^{-1}(j)$ and we get a cumulative memory bound for the last k segments of at least

$$(\sigma_{\ell-k+1} + \sum_{j=1}^{k-1} h_0^{-1}(j)) \cdot H. \quad (7)$$

▷ **Claim 4.9.** $\sigma_{\ell-k+1} + \sum_{j=1}^{k-1} h_0^{-1}(j) \geq \mathcal{L}_{h_0}(S^{max}) \cdot \sigma_{\ell-k+1} \cdot k$.

Proof of Claim. Observe that it suffices to prove the claim when we replace $\sigma_{\ell-k+1}$, which appears on both sides, by a larger quantity. In particular, we show how to prove the claim with $h_0^{-1}(k)$ instead, which is larger since $h_0(\sigma_{\ell-k+1}) < k$. But this follows immediately since by definition $\mathcal{L}_{h_0}(S^{max}) \leq \frac{\sum_{j=1}^k h_0^{-1}(j)}{k \cdot h_0^{-1}(k)}$, which is equivalent to what we want to prove. ◁

Write $S_{i^*} = \sigma_{\ell-k+1}$. By the claim, the cumulative memory contribution associated with the last block beginning at t_{i^*} is at least $\mathcal{L}_{h_0}(S^{max}) \cdot S_{i^*} \cdot h_0(S_{i^*})H$.

We repeat this in turn to find the time step for the beginning of the next block from the end, t_{i^*-1} . One small difference now is that there is a last partial segment of height at most H from the beginning of segment containing t_{i^*} to layer t_{i^*} . However, this only adds at most $h_1(n)/2$ to the length of the segment which still remains well within the height bound of $h(S_{i^*-1}, n) = h_0(S_{i^*-1})h_1(n)$ for Property (B) to apply.

Repeating this back to the beginning of the branching program we obtain a decomposition of the branching program into some number i^* of blocks, the i -th block beginning at time step t_i with 2^{S_i} nodes, height between $h_0(S_i)H$ and $h_0(S_i)H + H \leq 2h_0(S_i)H$, and with an associated cumulative memory contribution in the i -th block of $\geq \mathcal{L}_{h_0}(S^{max}) \cdot S_i \cdot h_0(S_i)H$. (This is correct even for the partial block starting at time $t_1 = 0$ since $S_1 = 0$.) Since we know that $i^* \leq \ell$, for convenience, we also define $S_i = 0$ for $i^* + 1 \leq i \leq \ell$. Then, by definition

$$CM(P) \geq \mathcal{L}_{h_0}(S^{max}) \cdot \left(\sum_{i=1}^{i^*} S_i \cdot h_0(S_i) \right) \cdot H = \mathcal{L}_{h_0}(S^{max}) \cdot \left(\sum_{i=1}^{\ell} S_i \cdot h_0(S_i) \right) \quad (8)$$

$$\text{and} \quad \sum_{i=1}^{\ell} h_0(S_i) \leq T/H. \quad (9)$$

As in the previous argument for the simple case, for $i \leq i^*$, we define the target k_i for the number of output values produced in each time block to be the smallest integer such that $CK^{-k_i} \leq 2^{-S_i} \alpha / (2T^{c+1})$. That is, $k_i = \lceil (S_i + \log_2(2CT^{c+1}/\alpha)) / \log_2 K \rceil$.

⁴ Since we are working backwards from the end of the branching program and we do not know how many segments are included in each block, we don't actually know this index until things stop with $t_1 = 0$

If $k_i > m'$ for some i , then $S_i \geq m' \cdot \log_2 K - \log_2(2CT^{c+1}/\alpha) \geq (m' \log_2 K)/2$ since m' is $\omega(\log n)$ and $1/\alpha$ and T are $n^{O(1)}$. Therefore $h_0(S_i) \geq h'(m'/2, n)$ and hence

$$CM(P) \geq \frac{1}{2} \mathcal{L}_{h_0}(S^{max}) \cdot m' \cdot h'(m'/2, n) \cdot \log_2 K$$

Suppose instead that $k_i \leq m'$ for all $i \leq i^*$. Then, for $x \sim \mu$, for each $i \in [i^*]$ and each sub-branching program B rooted at some node in L_{t_i} and extending until time t_{i+1} , by our choice of k_i and Property (B), the probability that B produces at least k_i correct outputs on input x is at most $\alpha \cdot 2^{-S_i}/(2T^{c+1})$. Therefore, by a union bound, for $x \sim \mu$ the probability that P produces at least k_i correct outputs in the i -th time block on input x is at most

$$|L_{t_i}| \cdot \alpha \cdot 2^{-S_i}/(2T^{c+1}) = \alpha/(2T^{c+1})$$

and hence the probability for $x \sim \mu$ that there is some i such that P produces at least k_i correct outputs on input x during the i -th block is at most $\ell \cdot \alpha/(2T^{c+1}) < \alpha/(2T^c)$. Therefore, the probability for $x \sim \mu$ that P produces at most $\sum_{i=1}^{\ell} (k_i - 1)$ correct outputs in total on input x is $> 1 - \alpha/(2T^c)$.

Since, by Property (A) and the maximum error it allows, P must produce at least m correct outputs with probability at least $\alpha - \epsilon \geq \alpha - \alpha(1 - 1/(2T^c)) = \alpha/(2T^c)$ for $x \sim \mu$, we must have $\sum_{i=1}^{i^*} (k_i - 1) \geq m$. Using our definition of k_i we obtain

$$\sum_{i=1}^{i^*} (S_i + \log_2(2CT^{c+1}/\alpha)) \geq m \log_2 K.$$

This is the one place in the proof where there is a distinction between an arbitrary branching program and one that comes from a random access machine.

We first start with the case of arbitrary branching programs: Note that $i^* \leq \ell = \lceil T/H \rceil = \lceil T/\lfloor h_1(n)/2 \rfloor \rceil$. Suppose that $T \log_2(2CT^{c+1}/\alpha) \leq \frac{1}{6} m \cdot h_1(n) \cdot \log_2 K$. Then, even with rounding, we obtain $\sum_{i=1}^{i^*} S_i \geq \frac{1}{2} m \log_2 K$.

Unlike an arbitrary branching program that may do non-trivial computation with sub-logarithmic S_i , a random-access machine with even one register requires at least $\log_2 n$ bits of memory (just to index the input for example) and hence $S_i + \log_2(2CT^{c+1}/\alpha)$ will be $O(S_i)$, since T is at most polynomial in n and $1/\alpha$ is at most polynomial in n by assumption. Therefore we obtain that $\sum_{i=1}^{i^*} S_i$ is $\Omega(m \log_2 K)$ without the assumption on T .

In the remainder we continue the argument for the case of arbitrary branching programs and track the constants involved. The same argument obviously applies for programs coming from random-access machines with slightly different constants that we will not track. In particular, since $S_i = 0$ for $i > i^*$ we have

$$\sum_{i=1}^{\ell} S_i \geq \frac{1}{2} m \cdot \log_2 K. \tag{10}$$

From this point we need to do something different from the argument in the simple case because the lower bound on the total cumulative memory contribution is given by Equation (8) and is not simply $\sum_{i=1}^{\ell} S_i \cdot H$. Instead, we combine Equation (10) and Equation (9) using the following technical lemma that we prove in Appendix A.

► **Lemma 4.10.** *Let $p : \mathbb{R}^{\geq 0} \rightarrow \mathbb{R}^{\geq 0}$ be a differentiable function such that $q(x) = x/p(x)$ is a concave increasing function of x . For $x_1, x_2, \dots \in \mathbb{R}^{\geq 0}$, if $\sum_i x_i \geq K$ and $\sum_i p(x_i) \leq L$ then $\sum_i x_i p(x_i) \geq q^{-1}(K/L) \cdot L$.*

17:16 Computational Cumulative Memory Lower Bounds

In our application of Lemma 4.10, $p = h_0$, $K = \frac{1}{2} m \cdot \log_2 K$, and $L = T/H$. Let S^* be the solution to $\frac{S^*}{h_0(S^*)} = K/L = \frac{m \cdot H \cdot \log_2 K}{2T} \geq \frac{m \cdot h_1(n) \log_2 K}{6T}$. Then Lemma 4.10 implies that $\sum_{i=1}^{\ell} S_i \cdot h_0(S_i) \geq S^* \cdot T/H = \frac{1}{2} m \cdot h_0(S^*) \cdot \log_2 K$. and hence

$$CM(P) \geq \mathcal{L}_{h_0}(S^{max}) \cdot \frac{1}{2} m \cdot h_0(S^*) \cdot H \cdot \log_2 K \geq \frac{1}{6} \mathcal{L}_{h_0}(S^{max}) \cdot m \cdot h(S^*, n) \cdot \log_2 K$$

since $H = \lfloor h_1(n)/2 \rfloor$ and $h(S^*, n) = h_0(S^*) \cdot h_1(n)$. \blacktriangleleft

In the special case that $h_0(s) = s^\Delta$ (and indeed for any nice function h_0), there is an alternative variant of the above in which one breaks up time into exponentially growing segments starting with time step T . We used that alternative approach in Section 3.

► **Remark 4.11.** If we restrict our attention to $o(m' \log K)$ -space bounded computation, then each $k_i \leq m'$ and the cumulative memory bound for a branching program in Theorem 4.8 becomes $\frac{1}{6} \mathcal{L}_{h_0}(n \log_2 |D|) \cdot m \cdot h(S^*(T, n), n) \cdot \log_2 K$. And the bound for RAM cumulative memory becomes $\Omega(\mathcal{L}_{h_0}(n \log_2 |D|) \cdot m \cdot h(S^*(T, n), n) \cdot \log_2 K)$.

Generic method for quantum time-space tradeoffs

Quantum circuit time-space lower bounds have the same general structure as their classical branching program counterparts. They require a lemma similar to (B) that gives an exponentially small probability of producing k outputs with a small number of queries.

► **Lemma 4.12** (Quantum generic property). *For all $k \leq m'$ and any quantum circuit \mathcal{C} with at most $h'(k, n)$ layers, there exists a distribution μ such that when $x \sim \mu$, the probability that \mathcal{C} produces at least k correct output values of $f(x)$ is at most $C \cdot K^{-k}$.*

Such lemmas have historically been proving using direct product theorems [29, 13] or the recording query technique [27]. Quantum time-space tradeoffs use the same blocking strategy as branching programs; however, they cannot use union bounds to account for input-dependent state at the start of a block. Instead, Proposition 3.4 lets us apply Lemma 4.12 to blocks in the middle of a quantum circuit. The 2^{2S} factor in Proposition 3.4 means that a quantum time-space or cumulative memory lower bound is half of what you would expect from a classical bound. Since a quantum circuit requires $\log_2 n$ qubits to make a query, we know that the space between layers is always at least $\log_2 n$ and

$$T \cdot S \text{ is } \Omega(\min\{m h'(S, n), m' Q(f)\} \cdot \log_2 K)$$

where $Q(f)$ is the bounded-error quantum query complexity of f .

Generic method for quantum cumulative complexity bounds

Our generic argument can just as easily be applied to quantum lower bounds for problems where we have an instance of Lemma 4.12 using Proposition 3.4 to bound the number of outputs produced even with initial input-dependent state. Quantum circuits require at least $\log_2 n$ qubits to hold the query index so the bounds derived are like those from Theorem 4.8(b).

► **Corollary 4.13.** *Let $c > 0$. Suppose that function f defined on D^n satisfies generic Lemma 4.12 with m' that is $\omega(\log_2 n)$. For $s > 0$, let $h(s, n) = h'(s/\log_2 K, n)$. Let $h(s, n) = h_0(s)h_1(n)$ where $h_0(1) = 1$ and h_0 is constant or a differentiable function where $s/h_0(s)$ is increasing and concave. Let S^* be defined by:*

$$S^*/h_0(S^*) = (m h_1(n) \log_2 K) / 6T.$$

The CMC used by a quantum circuit that computes f in time T with error $\varepsilon \leq (1 - 1/(2T^c))$ is

$$\geq \frac{1}{6} \mathcal{L}_{h_0}(n \log_2 |D|) \cdot \min \{m h(S^*, n), 3m' h'(m'/2, n)\} \cdot \log_2 K.$$

If the circuit uses $o(m' \log K)$ qubits, the CMC instead is $\frac{1}{6} \mathcal{L}_{h_0}(n \log_2 |D|) \cdot m \cdot h(S^*, n) \cdot \log_2 K$.

Full general applications. Some applications of our full general theorem generalizing classical and quantum time-space product lower bounds are the following (full proofs are in [18]):

► **Corollary 4.14.** *The CMC for listing all uniquely occurring elements or sorting n integers from $[n]$ is $\Omega(n^2)$.*

► **Corollary 4.15.** *Matrix multiplication for $n \times n$ matrices over finite field \mathbb{F} requires classical CMC that is $\Omega((n^6 \log |\mathbb{F}|)/T)$.*

► **Corollary 4.16.** *For every $m \geq n$, finding k disjoint collisions in a random function from $[m]$ to $[n]$ requires quantum CMC that is $\Omega(k^3 n/T^2)$.*

Corollary 4.14 uses properties (A) and (B) for unique elements with $h'(k, n) = n/4$, $m' = n/4$, $m = n/(2e)$, $\alpha = 1/(2e - 1)$, $K = 1/(2 \ln N)$ and $C = 1$ that follow from [16, Lemmas 2, 3]. Corollary 4.15 uses properties (A) and (B') with $h'(k, n) = \Theta(n\sqrt{k})$, $m' = m = n^2$, $\alpha = 1$, $K = |\mathbb{F}|^{\Theta(1)}$ and $C = d^2$. as proven in [4, Theorem 7.1]. Corollary 4.16 uses Lemma 4.12 with $h'(k, n) = \Theta(k^{2/3}n^{1/3})$, $m' = m = k$, and $C = K = 2$ which follow from [27, Theorem 9].

References

- 1 Scott Aaronson. Limitations of quantum advice and one-way communication. *Theory of Computing*, 1(1):1–28, 2005. doi:10.4086/toc.2005.v001a001.
- 2 Karl R. Abrahamson. Generalized string matching. *SIAM J. Comput.*, 16(6):1039–1051, 1987. doi:10.1137/0216067.
- 3 Karl R. Abrahamson. A time-space tradeoff for Boolean matrix multiplication. In *31st Annual IEEE Symposium on Foundations of Computer Science, Volume I*, pages 412–419, 1990. doi:10.1109/FSCS.1990.89561.
- 4 Karl R. Abrahamson. Time-space tradeoffs for algebraic problems on general sequential machines. *J. Comput. Syst. Sci.*, 43(2):269–289, 1991. doi:10.1016/0022-0000(91)90014-v.
- 5 Miklós Ajtai. Determinism versus nondeterminism for linear time RAMs with memory restrictions. *J. Comput. Syst. Sci.*, 65(1):2–37, 2002. doi:10.1006/jcss.2002.1821.
- 6 Miklós Ajtai. A non-linear time lower bound for Boolean branching programs. *Theory Comput.*, 1(1):149–176, 2005. doi:10.4086/toc.2005.v001a008.
- 7 Joël Alwen and Jeremiah Blocki. Efficiently computing data-independent memory-hard functions. In *Advances in Cryptology – CRYPTO 2016*, pages 241–271, 2016.
- 8 Joël Alwen, Jeremiah Blocki, and Krzysztof Pietrzak. Depth-robust graphs and their cumulative memory complexity. In *Advances in Cryptology – EUROCRYPT 2017*, pages 3–32, 2017.
- 9 Joël Alwen, Binyi Chen, Chethan Kamath, Vladimir Kolmogorov, Krzysztof Pietrzak, and Stefano Tessaro. On the complexity of Scrypt and proofs of space in the parallel random oracle model. In *Advances in Cryptology - EUROCRYPT 2016, Proceedings, Part II*, volume 9666 of *LNCS*, pages 358–387, 2016. doi:10.1007/978-3-662-49896-5_13.
- 10 Joël Alwen, Binyi Chen, Krzysztof Pietrzak, Leonid Reyzin, and Stefano Tessaro. Scrypt is maximally memory-hard. In *Advances in Cryptology - EUROCRYPT 2017, Proceedings, Part III*, volume 10212 of *Lecture Notes in Computer Science*, pages 33–62, 2017. doi:10.1007/978-3-319-56617-7_2.

- 11 Joël Alwen, Susanna F. de Rezende, Jakob Nordström, and Marc Vinyals. Cumulative space in black-white pebbling and resolution. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*, volume 67 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 38:1–38:21, 2017. doi:10.4230/LIPIcs.ITCS.2017.38.
- 12 Joël Alwen and Vladimir Serbinnenko. High parallel complexity graphs and memory-hard functions. In *Proceedings of the Forty-Seventh Annual ACM Symposium on Theory of Computing*, pages 595–603, 2015. doi:10.1145/2746539.2746622.
- 13 Andris Ambainis, Robert Špalek, and Ronald de Wolf. A new quantum lower bound method, with applications to direct product theorems and time-space tradeoffs. *Algorithmica*, 55(3):422–461, 2009. doi:10.1007/s00453-007-9022-9.
- 14 Mohammad Hassan Ameri, Alexander R. Block, and Jeremiah Blocki. Memory-hard puzzles in the standard model with applications to memory-hard functions and resource-bounded locally decodable codes. Cryptology ePrint Archive, Paper 2021/801, 2021. URL: <https://eprint.iacr.org/2021/801>.
- 15 Andrew Baird, Bryant Bost, Stefano Buliani, Vyom Nagrani, Ajay Nair, Rahul Papat, and Brajendra Singh. AWS serverless multi-tier architectures with Amazon API Gateway and AWS Lambda, 2021. URL: <https://docs.aws.amazon.com/whitepapers/latest/serverless-multi-tier-architectures-api-gateway-lambda/welcome.html>.
- 16 Paul Beame. A general sequential time-space tradeoff for finding unique elements. *SIAM J. Comput.*, 20(2):270–277, 1991. doi:10.1137/0220017.
- 17 Paul Beame, T. S. Jayram, and Michael E. Saks. Time-space tradeoffs for branching programs. *J. Comput. Syst. Sci.*, 63(4):542–572, 2001. doi:10.1006/jcss.2001.1778.
- 18 Paul Beame and Niels Kornerup. Cumulative memory lower bounds for randomized and quantum computation. *CoRR*, abs/2301.05680, 2023. doi:10.48550/arXiv.2301.05680.
- 19 Paul Beame, Michael E. Saks, Xiaodong Sun, and Erik Vee. Time-space trade-off lower bounds for randomized computation of decision problems. *J. ACM*, 50(2):154–195, 2003. doi:10.1145/636865.636867.
- 20 Jeremiah Blocki and Samson Zhou. On the depth-robustness and cumulative pebbling cost of Argon2i. In *Theory of Cryptography*, pages 445–465, 2017.
- 21 Dan Boneh, Henry Corrigan-Gibbs, and Stuart Schechter. Balloon hashing: A memory-hard function providing provable protection against sequential attacks. In *Advances in Cryptology – ASIACRYPT 2016*, pages 220–248, 2016.
- 22 Allan Borodin and Stephen A. Cook. A time-space tradeoff for sorting on a general sequential model of computation. *SIAM J. Comput.*, 11(2):287–297, 1982. doi:10.1137/0211022.
- 23 Allan Borodin, Michael J. Fischer, David G. Kirkpatrick, Nancy A. Lynch, and Martin Tompa. A time-space tradeoff for sorting on non-oblivious machines. *J. Comput. Syst. Sci.*, 22(3):351–364, 1981. doi:10.1016/0022-0000(81)90037-4.
- 24 Allan Borodin, Alexander A. Razborov, and Roman Smolensky. On lower bounds for read-k-times branching programs. *Comput. Complex.*, 3:1–18, 1993. doi:10.1007/BF01200404.
- 25 Binyi Chen and Stefano Tessaro. Memory-hard functions from cryptographic primitives. In *Advances in Cryptology – CRYPTO 2019*, pages 543–572, 2019.
- 26 Cynthia Dwork, Moni Naor, and Hoeteck Wee. Pebbling and proofs of work. In *Advances in Cryptology – CRYPTO 2005*, pages 37–54, 2005.
- 27 Yassine Hamoudi and Frédéric Magniez. Quantum time-space tradeoff for finding multiple collision pairs. In *16th Conference on the Theory of Quantum Computation, Communication and Cryptography (TQC 2021)*, volume 197 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 1:1–1:21, 2021. doi:10.4230/LIPIcs.TQC.2021.1.
- 28 Stasys Jukna. A nondeterministic space-time tradeoff for linear codes. *Inf. Process. Lett.*, 109(5):286–289, 2009. doi:10.1016/j.ipl.2008.11.001.
- 29 Hartmut Klauck, Robert Špalek, and Ronald de Wolf. Quantum and classical strong direct product theorems and optimal time-space tradeoffs. *SIAM Journal on Computing*, 36(5):1472–1493, 2007. doi:10.1137/05063235x.

- 30 Yishay Mansour, Noam Nisan, and Prason Tiwari. The computational complexity of universal hashing. *Theor. Comput. Sci.*, 107(1):121–133, 1993. doi:10.1016/0304-3975(93)90257-T.
- 31 E. Okol'nishnikova. On lower bounds for branching programs. *Siberian Advances in Mathematics*, 3(1):152–166, 1993.
- 32 Ling Ren and Srinivas Devadas. Proof of space from stacked expanders. In *Proceedings, Part I, of the 14th International Conference on Theory of Cryptography - Volume 9985*, pages 262–285. Springer-Verlag, 2016. doi:10.1007/978-3-662-53641-4_11.
- 33 Martin Sauerhoff and Philipp Woelfel. Time-space tradeoff lower bounds for integer multiplication and graphs of arithmetic functions. In *Proceedings of the 35th Annual ACM Symposium on Theory of Computing*, pages 186–195, 2003. doi:10.1145/780542.780571.
- 34 John E. Savage. *Models of Computation: Exploring the Power of Computing*. Addison-Wesley Longman Publishing Co., Inc., USA, 1st edition, 1997.
- 35 Martin Tompa. Time-space tradeoffs for computing functions, using connectivity properties of their circuits. *J. Comput. Syst. Sci.*, 20(2):118–132, 1980. doi:10.1016/0022-0000(80)90056-2.
- 36 Andrew Chi-Chih Yao. Probabilistic computations: Toward a unified measure of complexity (extended abstract). In *18th Annual IEEE Symposium on Foundations of Computer Science*, pages 222–227, 1977. doi:10.1109/sfcs.1977.24.
- 37 Yaacov Yesha. Time-space tradeoffs for matrix multiplication and the discrete Fourier transform on any general sequential random-access computer. *Journal of Computer and System Sciences*, 29(2):183–197, 1984. doi:10.1016/0022-0000(84)90029-1.
- 38 Mark Zhandry. How to record quantum queries, and applications to quantum indifferentiability. In *Advances in Cryptology – CRYPTO 2019*, pages 239–268, 2019.

A Optimizations

In this section we prove general optimization lemmas that allow us to derive worst-case properties of the allocation of branching program layers into blocks.

► **Lemma A.1.** *For non-negative reals x_1, x_2, \dots if $\sum_i x_i \leq \sum_i x_i^2$ then $\sum_i x_i^3 \geq \sum_i x_i^2$.*

Proof. Without loss generality we remove all x_i that are 0 or 1 since they contribute the same amount to each of $\sum_i x_i$, $\sum_i x_i^2$, and $\sum_i x_i^3$. Therefore every x_i satisfies $0 < x_i < 1$ or it satisfies $x_i > 1$. We rename those x_i with $0 < x_i < 1$ by y_i and those x_i with $x_i > 1$ by z_j .

Then $\sum_i x_i \leq \sum_i x_i^2$ can be rewritten as $\sum_i y_i(1-y_i) \leq \sum_j z_j(z_j-1)$, and both quantities are positive. Let y^* be the largest value < 1 and z^* be the smallest value > 1 . Thus:

$$\begin{aligned} \sum_i (y_i^2 - y_i^3) &= \sum_i y_i^2(1 - y_i) \leq \sum_i y^* y_i(1 - y_i) = y^* \sum_i y_i(1 - y_i) \leq y^* \sum_j z_j(z_j - 1) \\ &< z^* \sum_j z_j(z_j - 1) = \sum_j z^* z_j(z_j - 1) \leq \sum_j z_j^2(z_j - 1) = \sum_j (z_j^3 - z_j^2). \end{aligned}$$

Rewriting gives $\sum_i y_i^2 + \sum_j z_j^2 < \sum_i y_i^3 + \sum_j z_j^3$, or $\sum_i x_i^3 > \sum_i x_i^2$, as required. ◀

The following is a generalization of the above to all differentiable functions $p : \mathbb{R}^{\geq 0} \rightarrow \mathbb{R}^{\geq 0}$ such that $s/p(s)$ is a concave increasing function of s .

► **Lemma 4.10.** *Let $p : \mathbb{R}^{\geq 0} \rightarrow \mathbb{R}^{\geq 0}$ be a differentiable function such that $q(x) = x/p(x)$ is a concave increasing function of x . For $x_1, x_2, \dots \in \mathbb{R}^{\geq 0}$, if $\sum_i x_i \geq K$ and $\sum_i p(x_i) \leq L$ then $\sum_i x_i p(x_i) \geq q^{-1}(K/L) \cdot L$.*

17:20 Computational Cumulative Memory Lower Bounds

Proof. By hypothesis, $\sum_i (x_i - Kp(x_i)/L) \geq 0$. Observe that $s - Kp(s)/L$ is an increasing function of s since $s/p(s)$ is an increasing function of s that is 0 precisely when $s = q^{-1}(K/L)$. Since all x_i with $x_i = q^{-1}(K/L)$ evaluate to 0 in the sum, we can rewrite it as

$$\sum_{x_i > q^{-1}(K/L)} (x_i - Kp(x_i)/L) \geq \sum_{x_i < q^{-1}(K/L)} (Kp(x_i)/L - x_i), \quad (11)$$

where each of the summed terms is positive. For $x_i \neq q^{-1}(K/L)$, define

$$f(x_i) = x_i \cdot \frac{p(x_i) - q^{-1}(K/L) \cdot L/K}{x_i - Kp(x_i)/L}.$$

Observe that for $x_i = q^{-1}(K/L)$ the denominator is 0 and the numerator equals $p(x_i) - x_i \cdot L/K$ which is also 0. For $x_i > q^{-1}(K/L)$ both the numerator and denominator are positive and for $x_i < q^{-1}(K/L)$ both the numerator and denominator are negative. Hence $f(x_i)$ is non-negative for every $x_i \neq q^{-1}(K/L)$. The following claim holds because of the concavity of q ; its proof is in the full paper [18].

▷ **Claim A.2.** If q is a convex differentiable function, we can complete f to a (non-decreasing) continuous function of x with $f'(x) \geq 0$ for all x with $0 < x \neq q^{-1}(K/L)$

We now have the tools we need. Let x_-^* be the largest $x_i < q^{-1}(K/L)$ and x_+^* be the smallest $x_i > q^{-1}(K/L)$. Then we have $f(x_+^*) \geq f(x_-^*)$ and

$$\begin{aligned} & \sum_{x_i > q^{-1}(K/L)} (x_i p(x_i) - q^{-1}(K/L) \cdot L/K \cdot x_i) \\ &= \sum_{x_i > q^{-1}(K/L)} f(x_i) \cdot (x_i - Kp(x_i)/L) \\ &\geq \sum_{x_i > q^{-1}(K/L)} f(x_+^*) \cdot (x_i - Kp(x_i)/L) \\ &\geq f(x_-^*) \sum_{x_i > q^{-1}(K/L)} (x_i - Kp(x_i)/L) \\ &\geq f(x_-^*) \sum_{x_i < q^{-1}(K/L)} (Kp(x_i)/L - x_i) \quad \text{by Equation (11)} \\ &\geq \sum_{x_i < q^{-1}(K/L)} f(x_i) \cdot (Kp(x_i)/L - x_i) \\ &= \sum_{x_i < q^{-1}(K/L)} (q^{-1}(K/L) \cdot L/K \cdot x_i - x_i p(x_i)). \end{aligned}$$

Adding back the terms where $x_i = q^{-1}(K/L)$, which have value 0, and rewriting we obtain

$$\sum_i (x_i p(x_i) - q^{-1}(K/L) \cdot L/K \cdot x_i) \geq 0.$$

Therefore we have

$$\sum_i x_i p(x_i) \geq q^{-1}(K/L) \cdot L/K \cdot \sum_i x_i \geq q^{-1}(K/L) \cdot (L/K) \cdot K = q^{-1}(K/L) \cdot L. \quad \blacktriangleleft$$