

# Zero-One Frequency Laws

Vladimir Braverman<sup>\*</sup>  
UCLA  
vova@cs.ucla.edu

Rafail Ostrovsky<sup>†</sup>  
UCLA  
rafail@cs.ucla.edu

## ABSTRACT

Data streams emerged as a critical model for multiple applications that handle vast amounts of data. One of the most influential and celebrated papers in streaming is the “AMS” paper on computing frequency moments by Alon, Matias and Szegedy. The main question left open (and explicitly asked) by AMS in 1996 is to give the precise characterization for which functions  $G$  on frequency vectors  $m_i$  ( $1 \leq i \leq n$ ) can  $\sum_{i \in [n]} G(m_i)$  be approximated efficiently, where “efficiently” means by a single pass over data stream and poly-logarithmic memory. No such characterization was known despite a tremendous amount of research on frequency-based functions in streaming literature. In this paper we finally resolve the AMS main question and give a precise characterization (in fact, a zero-one law) for *all* monotonically increasing functions on frequencies that are zero at the origin.

That is, we consider all monotonic functions  $G: R \mapsto R$  such that  $G(0) = 0$  and  $G$  can be computed in poly-logarithmic time and space and ask, for which  $G$  in this class is there an  $(1 \pm \epsilon)$ -approximation algorithm for computing  $\sum_{i \in [n]} G(m_i)$  for any polylogarithmic  $\epsilon$ ? We give an algebraic characterization for all such  $G$  so that:

- For all functions  $G$  in our class that satisfy our algebraic condition, we provide a very general and constructive way to derive an efficient  $(1 \pm \epsilon)$ -approximation algorithm for computing  $\sum_{i \in [n]} G(m_i)$  with polylogarithmic memory and a single pass over data stream; while
- For all functions  $G$  in our class that *do not* satisfy our algebraic characterization, we show a lower bound

<sup>\*</sup>Supported in part by NSF grants 0716835, 0716389, 0830803, 0916574.

<sup>†</sup>Supported in part by IBM Faculty Award, Xerox Innovation Group Award, the Okawa Foundation Award, Lockheed Martin, Intel, Teradata, NSF grants 0716835, 0716389, 0830803, 0916574, BSF grant and U.C. MICRO grant.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

STOC'10, June 5–8, 2010, Cambridge, Massachusetts, USA.  
Copyright 2010 ACM 978-1-4503-0050-6/10/06 ...\$10.00.

that requires greater than polylog memory for computing an approximation to  $\sum_{i \in [n]} G(m_i)$  by *any* one-pass streaming algorithm.

Thus, we provide a zero-one law for all monotonically increasing functions  $G$  which are zero at the origin. Our results are quite general. As just one illustrative example, our main theorem implies a lower bound for  $G(x) = (x(x-1))^{0.5 \arctan(x+1)}$ , while for a function  $G(x) = (x(x+1))^{0.5 \arctan(x+1)}$  our main theorem automatically yields a polylog memory one-pass  $(1 \pm \epsilon)$ -approximation algorithm for computing  $\sum_{i \in [n]} G(m_i)$ . For both of these examples no lower or upper bounds were known. Of course, these are just illustrative examples, and there are many others. One might argue that these two functions may not be of interest in practical applications – we stress that our law works for **all** functions in this class, and the above examples illustrate the power of our method.

To the best of our knowledge, this is the first zero-one law in the streaming model for a wide class of functions, though we suspect that there are many more such laws to be discovered. Surprisingly, our upper bound requires only 4-wise independence and does not need the stronger machinery of Nisan’s pseudorandom generators, even though our class captures multiple functions that previously required Nisan’s generators. Furthermore, we believe that our methods can be extended to the more general models and complexity classes. For instance, the law also holds for a smaller class of non-decreasing and symmetric functions (i.e.,  $G(x) = G(-x)$  and  $G(0) = 0$ ) which, due to negative values, allow deletions.

## Categories and Subject Descriptors

F.2 [Theory of Computation]: ANALYSIS OF ALGORITHMS AND PROBLEM COMPLEXITY

## General Terms

Algorithms, Theory.

## Keywords

Data Streams, Randomized Algorithms, Theory of Computation.

# 1. INTRODUCTION

Data streams emerged as a critical model for many applications with vast amounts of data. The importance of the streaming model is discussed, e.g., by Aggarwal (ed.) [1] and Muthukrishnan [41]. In the seminal AMS paper, Alon, Matias and Szegedy [2] studied the following basic model:

**DEFINITION 1.1.** *Let  $m, n$  be positive integers. A stream  $D = D(n, m)$  is a sequence of size  $m$  of integers  $p_1, \dots, p_m$ , where  $p_i \in \{1, \dots, n\}$ . A frequency vector  $M = M(D)$  is a vector of dimensionality  $n$  with non-negative entries  $m_i, i \in [n]$  defined as:*

$$m_i = |\{j : 1 \leq j \leq m, p_j = i\}|.$$

AMS [2] studied the problem of approximating frequency moments with sublinear memory. They showed<sup>1</sup> that for  $k = 0, 1, 2$  it is possible to approximate  $F_k$  with polylogarithmic space; for  $k > 2$  they gave  $O^*(n^{1-1/k})$  upper bound. Also, they gave  $O^*(n^{1-5/k})$  lower bound for any  $k > 5$ . Indyk [30] presented a celebrated method of stable distributions for approximating  $L_p$  norms  $p \in (0, 2]$  in a general model where deletions are allowed and updates can be larger than 1. Indyk and Woodruff [33] gave the first optimal algorithm for  $F_k, k > 2$ , proving  $O^*(n^{1-2/k})$  upper bound. This result was later improved by polylog factors by Bhuvanagiri, Ganguly, Kesh and Saha [7]. Bar-Yossef, Jayram, Kumar and Sivakumar [3] used information theory to prove the first nearly matching lower bound of  $\Omega(n^{1-(2+\epsilon)/k})$ . Later Chakrabarti, Khot and Sun [17] improved the lower bound to  $\Omega(n^{1-2/k})$  for one-pass algorithms. Indyk and Woodruff [32] and Woodruff [43] gave optimal lower bound in terms of error parameter. Many other results on frequency moments include, e.g., Flajolet and Martin [24], Bar-Yossef, Jayram, Kumar, and Sivakumar [4], Coppersmith and Kumar [19], Cormode, Datar, Indyk and Muthukrishnan [20], Feigenbaum, Kannan, Strauss and Viswanathan [23], Ganguly [25], Ganguly and Cormode [26], Li [39], and Kane, Nelson and Woodruff [35, 36], Braverman and Ostrovsky [9, 11]. Currently, many important frequency-based functions are well-understood in different models. Research on entropy, entropy norm and distributions included the works of Bhuvanagiri and Ganguly [6], Chakrabarti, Do Ba and Muthukrishnan [14], Chakrabarti, Cormode and McGregor [15], Guha, McGregor and Venkatasubramanian [28], Harvey, Nelson and Onak [29], Indyk and McGregor [31], Lall, Sekar, Ogihara, Xu and Zhang [38], Braverman, Chung, Liu, Mitzenmacher and Ostrovsky [8], Braverman and Ostrovsky [13]. The related question of frequent elements has been studied by Charikar, Chen and Farach-Colton [18], Cormode and Hadjieleftheriou [21], Cormode and Muthukrishnan [22]. The frequency-based functions were studied in extended models such as the read/write model (Beame, Jayram and Rudra [5]), and the randomized model (Chakrabarti, Cormode and McGregor [16], Jayram, McGregor, Muthukrishnan and Vee [34]).

The main question left open (and explicitly asked) by Alon, Matias and Szegedy [2] is:

**AMS (informal):** *What other frequency-based functions can be approximated on streams?*

<sup>1</sup>This is a very informal explanation. For precise statements see [2].

In 2006 Guha and Indyk [40] (Question 5) asked a related question:

**Guha, Indyk (informal):** *What distances can be computed between two distribution vectors  $V$  and  $U$  defined by two streams? Consider a function  $\phi(x, y)$  such that  $\phi(x, x) = 0$ . What are the properties of  $\phi$  such that  $\sum_{i \in [n]} \phi(v_i, u_i)$  can be approximated with small memory and multiplicative error?*

To the best of our knowledge, the only work in this direction is the result of Guha, Indyk and McGregor [27]. They proved the Shift Invariant Theorem, a general result that gives a necessary condition for approximating a wide class of two-variable functions. The Shift Invariant Theorem says, very informally, that if  $\frac{\phi(a, a+c)}{\phi(b, b+c)} > \alpha > 1$  for some  $a, b, c$ , then a linear lower bound can be shown. In fact, they showed a stronger result for *distributions*, i.e., vectors with entries  $m_i/m$ . As a corollary, they obtained a linear lower bound for all commonly used *divergences* and gave additive approximations instead. Guha, Indyk and McGregor suggested that any function  $\phi(x, y)$  is not sketchable unless it is a function of  $x-y$ . Thus, it is important to understand what functions of a single variable can be approximated. Unfortunately, for such functions, the Shift Invariant Theorem is not applicable.

The questions of Alon, Matias and Szegedy [2] and Guha and Indyk [40] are still unresolved. Indeed, no sufficient and necessary conditions are known for a function to be computable; not even in the basic model of Alon, Matias and Szegedy [2]. Thus, to even approach the question posed by Guha and Indyk, we must first resolve the question of AMS [2]. This is exactly what we do in this paper. We limit ourselves to the basic streaming model (Definition 1.1) and consider a class of non-decreasing functions  $G : R \mapsto R$  such that  $G(0) = 0$ . We discuss the following problem:

G-sum( $D, \epsilon$ )

Given: A stream  $D = D(n, m)$  and  $\epsilon = \Omega^*(1)^a$ .

Output: A real  $r$  such that  $r$  is a  $(1 \pm \epsilon)$ -approximation<sup>b</sup> of  $\sum_{i \in [n]} G(m_i)$ .

<sup>a</sup>Given the parameters  $n, m$ , we use  $O_{n,m}^*(g(n, m))$  to denote  $O(\log^{O(1)}(nm)g(n, m))$ ;  $\Omega_{n,m}^*(g(n, m))$  is defined similarly. When the context is clear, we use  $O^*$  instead of  $O_{n,m}^*$ .

<sup>b</sup>We say that  $x$  is a  $(1 \pm \epsilon)$ -approximation of  $y$  if  $(1 - \epsilon)y \leq x \leq (1 + \epsilon)y$ .

We ask whether or not it is possible to solve the G-Sum problem in polylog space, with a single pass and a small probability of error. More precisely, we define the following class **STREAM-POLYLOG** class.

**DEFINITION 1.2.** *We say that a function  $G$  belongs to the **STREAM-POLYLOG** class if for any  $k$  there exist  $t$  and an algorithm  $\mathfrak{A}$  such that for any  $n, m$ ; any stream  $D = D(n, m)$ ; and any<sup>2</sup>  $\epsilon \geq \frac{1}{\log^k(nm)}$ :*

<sup>2</sup>Here, polylog  $\epsilon$  can also be replaced with a larger error which gives additional results as well; we defer these generalizations to the expanded version.

1.  $\mathfrak{A}$  makes a single pass over  $D$ .
2.  $\mathfrak{A}$  solves  $G$ -Sum( $D, \epsilon$ ) and errs with probability at most<sup>3</sup> 0.3.
3.  $\mathfrak{A}$  uses  $O(\log^t(nm))$  memory bits.

Our main result is a sufficient and necessary condition (i.e., zero-one law) for any non-decreasing  $G$  such that  $G(0) = 0$ . An important strength of our result is that we consider functions which may not even be differentiable or continuous and can include “jumps”:

DEFINITION 1.3. **Local jump.** *Define:*

$$\pi_\epsilon(x) = \min\{x, \min\{|z| \in \mathbb{N}^+ : |G(x) - G(x+z)| > \epsilon G(x)\}\}$$

Now, we are ready to introduce the following notion of “tractable” function:

DEFINITION 1.4. *Function  $G$  is Tractable if  $G(1) > 0$  and:*

$$\forall k \exists N_0 \exists t \forall x, y \in \mathbb{N}^+, \forall R \in \mathbb{R}^+ \forall \epsilon : \left( R > N_0, \frac{G(x)}{G(y)} = R, \epsilon > \frac{1}{\log^k(Rx)} \right) \implies \left( \left( \frac{\pi_\epsilon(x)}{y} \right)^2 \geq \frac{R}{\log^t(xR)} \right). \quad (1)$$

Our main result is the following theorem:

THEOREM 1.5. **Main Theorem** *Let  $G$  be a non-decreasing function such that  $G(0) = 0$ :*

$G \in \text{STREAM-POLYLOG}$  if and only if  $G \in \text{Tractable}$

The precise formula of a tractable function is somewhat complicated; however, the notion of tractability is very intuitive, as we explain below.

HIGH-LEVEL INTERPRETATION OF OUR ZERO-ONE LAW: (We stress that here we only explain how our result should be interpreted, without getting into any technical details regarding the proof. We explain the intuition about our technical proof, that turned out to be highly nontrivial, in a different section below and explain why our result gives us the intuitively “right” characterization.) *The space complexity of  $G$ -Sum should depend on how fast  $G$  grows. Quadratic growth seems to be a “waterline” since any  $G(x) = x^{2+k}$  can be approximated if and only if  $k \leq 0$ . Can we prove that  $G(x)$  can be approximated if and only if  $G(x)$  grows “slower” than  $x^2$ ? A careful elaboration of the above question gives us (1). We need to define precisely the notion of “growing slower than.” It turns out that the right measure is the ratio between  $G(x)/G(y)$  and  $(x/y)^2$  for all pairs  $x, y$ . First, if  $G(1) = 0$  then the maximal ratio is infinity. Indeed, we can generate vectors of arbitrary large dimensionality and still preserve a constant value of  $G$ -Sum. Intuitively it is clear that such a function cannot be approximated in a small space. This is indeed the case; similar observations were made, for example, for entropy norm  $G(x) = x \log(x)$  in [14]. If  $G(1) > 0$ , then we can consider the ratio  $G(x)/G(1)$ . If it grows significantly faster (i.e., beyond polylog factors) than  $x^2$  then*

<sup>3</sup>We remark that 0.3 is for concreteness only and can be replaced by any  $1/\text{poly}$  probability.

we cannot approximate  $G$ -sum in polylog space (the lower bound follows from [17]). The above rule should be extended to any pair of positive integers  $x > y$ :  $G(x)/G(y)$  cannot be significantly and asymptotically larger than  $x^2/y^2$ . Finally, can  $G$  grow very fast locally? Consider the following step function:  $G(x) = 2^{\lfloor \log(x) \rfloor}$ . It is not hard to see that we cannot obtain a small constant approximation of this function with sublinear space. Thus, in addition to the “global” growth of  $G$ , we also need to address its “local” jumps. We take one step further and define  $\pi_\epsilon(x)$  as a minimal shift (to the left or the right of  $x$ ) needed to “jump” beyond the  $(1 \pm \epsilon)$  factor. It turns out that the ratio between  $G(x)/G(y)$  and  $\pi(x)^2/y^2$  is what we need to measure; this is exactly what (1) states.<sup>4</sup>

Our main contribution is the necessary and sufficient condition for a class of non-decreasing non-negative functions. To the best of our knowledge, this is the first zero-one law in the streaming model for a wide class of functions. In addition, our method has the following advantages.

RAMIFICATIONS OF OUR MAIN RESULT: Streaming algorithms typically employ a special treatment for every new function. For example, there are unique algorithms for  $F_2$  [2, 30],  $F_0$  [24, 20, 2] and for  $L_p, 0 \leq p < 2$  [30, 39, 35]. Every new function so far required a new method or modification of existing methods and proofs. There is a natural question that one can ask:

*Is there a general algorithm for all tractable functions?*

We take the first step towards answering this fascinating general question. Our main technical contribution is a general algorithm and general proof for any function that satisfies (1).

Many of the existing methods first assume that totally random vectors are available, and later employ the celebrated pseudorandom generator of Nisan [42] to reduce the space for randomness. This brings another natural question:

*Are pseudorandom generators necessary or can we directly work with  $k$ -wise independent distributions?*

Surprisingly, we show that our general algorithm requires only 4-wise independence.

GENERALIZATIONS: Our methods can be extended to more general models and complexity classes. Our main theorem is immediately applicable to a more general model with deletions and larger updates. Here the law holds for a smaller class of non-decreasing and symmetric functions (i.e.,  $G(x) = G(-x)$ ).

Our methods may also shed some light on the fascinating open problem posed by Guha and Indyk in 2006. In fact, we believe that our methods are quite general and extend to other classes such as sublinear-space functions or distributions (i.e., to the functions  $G(m_i/m)$ ).

RELATED WORK: Our paper is a generalization of many previous works and explicitly employs existing ideas as well as the development of several new general techniques, both for our algorithm and analysis. The major difference with

<sup>4</sup>In fact, the Shift Invariant Theorem gives a similar intuition for functions of two variables by measuring the ratio  $G(x, x+a)/G(y, y+a)$ .

previous works is that assumptions that may be true for specific functions are not true in general. This makes our task highly non-trivial. As a result, we needed to develop a general and novel framework with minimum assumptions. At the end, the only assumption we make is that  $G$  is non-decreasing.

In particular, our method was inspired by the seminal work of Indyk and Woodruff [33]; our work can be seen as a generalization and simplification of [33]. The most notable difference is that the Indyk and Woodruff’s method was specifically developed to tackle  $F_k$ . They employ the strong relation between  $m_i$  and  $m_i^k$ ; we cannot make this assumption. One example of Indyk and Woodruff method is a statement that if  $y$  is an  $(1 \pm \epsilon)$ -approximation of  $x$  then  $G(y)$  is a  $(1 \pm \epsilon')$ -approximation of  $G(x)$  for some  $\epsilon' = \epsilon^{O(1)}$ . The statement is correct for  $G(x) = x^k$  for any constant  $k$ ; but notice that it is not correct in general. In addition, and perhaps somewhat surprisingly, our method requires only 4-wise independence whereas as Indyk and Woodruff [33] employs pseudorandom generators of Nisan [42].

One of the key steps in our new method is an approximation of “heavy” elements  $G(m_i)$ . Heavy elements have been intensively studied, e.g., in the papers of Charikar, Chen and Farach-Colton [18] and Cormode and Muthukrishnan [22]. Unfortunately, these methods are not directly applicable: we need  $G$ -heaviness, where [18, 22] address  $L_2$  and  $L_1$  heaviness respectfully. We thus developed a novel method for  $G$ -heaviness. Our method uses ideas developed in [18] and new techniques; both our method and [18] employ and build upon the AMS sketching for  $F_2$  [2].

**MAIN TECHNICAL IDEAS:** The lower bound follows from the reduction to set disjointness and index problems. Assume that  $G(x)/G(y)$  is both arbitrarily large and significantly larger than  $(x/y)^2$ . Then set disjointness can be translated into a stream that contains only elements with frequency  $y$  or (in addition) one element with a frequency larger than  $x$ . Since  $G$  is non-decreasing, and with some additional work, by utilizing [17] we derive our lower bound. The detailed proof is technically more involved than the direct application of [17] and can be found in Section 3.

To prove the upper bound we apply two key steps: first, we reduce the problem of  $G$ -Sum to the question of estimating heavy elements; and second, we solve the heavy elements problem separately.

To address the first reduction, we generalize the machinery of Indyk and Woodruff. We split the domain of  $G$  into “intervals” and reduce  $G$ -Sum to the problem of counting the number of entries  $G(m_i)$  that belong to a single interval. In fact, only a relaxed variant of interval counting should be solved. We always need to maintain a multiplicative upper bound; however, the lower bound is needed only for intervals that contribute nontrivially to the final answer. This task can be solved (if the number of entries in the interval is small) by reduction to finding heavy elements. Finally, we show how to reduce general counting to counting on intervals with a small number of elements.

To complete the proof, we need to approximate any  $G(m_i)$  such that  $G(m_i) = \Omega^*(\sum_{j \neq i} G(m_j))$ . The definition of  $\pi_\epsilon(m_i)$  implies that it is sufficient to approximate  $m_i$  with additive error smaller than  $\pi_\epsilon(m_i)$ . We need to verify that  $G(m_i) = \Omega^*(\sum_{j \neq i} G(m_j))$ , and that the additive error is less than  $\pi_\epsilon(m_i)$ . These generally hard problems can be

solved for tractable functions. The key observation is that if  $G$  is tractable, then  $G(m_i) = \Omega^*(\sum_{j \neq i} G(m_j))$  implies that  $(\pi_\epsilon(m_i))^2 = \Omega^*(\sum_{j \neq i} m_j^2)$ . This implies that to find heavy elements we only need to verify that  $m_i^2 = \Omega^*(\sum_{j \neq i} m_j^2)$  and approximate  $m_i$  with additive error. With non-trivial technical effort we obtain an algorithm that reports only  $(1 \pm \epsilon)$ -approximations of  $G(m_i)$ , and with high probability finds all heavy  $G(m_i)$ . There is a possibility that the algorithm will report some extra  $G(m_i)$ . However, we show that the false positives do not affect the correctness of the final algorithm. Section 5 is devoted to this step.

## 2. PURIFEIR: AN ALTERNATIVE VIEW

In this paper and in our other STOC paper [13] we solve two very different problems. The models are incomparable: this paper addresses implicitly defined tensors, and [13] describes how to approximate functions on explicit streaming frequencies. Moreover, our objectives are different. In [13] we improve the existing  $\log(N)$ -approximation to an  $(1 + \epsilon)$ -approximation. In this paper we give new algorithms for previously unknown functions. As a result, our bottom-line techniques are different. Yet we use a single methodology for both problems, which seems quite general, and worth explaining.

Indyk and Woodruff proposed a seminal method of layering that has become folklore in the streaming applications. They show that the frequency moments problem can be reduced to a question of finding heavy hitters. Further, they show that it is possible to use  $F_2$  (in particular the algorithm of [18]) to compute heavy hitters under  $F_k$ . Although the Indyk and Woodruff algorithm solves a particular problem, we believe that their methodology goes far beyond frequency moments. In fact our two papers can be seen as building upon methodology of Indyk and Woodruff and in fact generalization of their methodology. Thus, we would like to emphasize the importance of their methodology and explain our contribution.

In many cases the most non-trivial step of streaming computations can be seen as a summation (or counting) of entries of a very large vector that is implicitly defined by the stream. We claim that the Indyk and Woodruff methodology is applicable to many such functions. We were able to show that the layering method can be used for at least two general settings: frequency-based functions and tensors. In fact, we believe that it is possible to specify the set of problems and conditions for which the layering reduction is possible, something which we are currently exploring.

It seems that one such condition is separability, i.e. an ability to efficiently sample the target vector (not the stream). Generally speaking, the entries of this imaginary vector may have a very complex nature and be dependent on each other. Thus, sampling of the stream does not necessarily correspond to the sampling of the implicit vector. Yet we show that for the Independence Problem from [13] it is possible to sample the vector without sampling the stream. Understanding which implicit vectors can be sampled efficiently is another important problem.

While the heavy hitters problem seems to be somewhat easier than to solve approximating the entire vector, it is generally not clear how to address this problem without computing the sum. Indeed, when the heavy hitter contributes almost entirely to the sum, those two problems co-

incide. Thus, it seems that in many cases finding the heavy hitter is the core obstacle. We are not aware of a general approach to tackling the heavy elements problem.<sup>5</sup> Thus, the key contribution of our methodology is such a method for tackling heavy elements. We present it very informally here; two specific examples of applying this methodology are in our two papers, and as these papers show, this methodology is quite flexible. Call an entry of a vector  $\alpha$ -significant if it contributes at least  $\alpha$ -fraction of the entire sum. Call a zero-one streaming function  $(\alpha, \beta)$ -certificate for the vector  $V = V(D)$  if  $\text{Certificate}(D) = 1$  implies that there is a  $\beta$ -significant element in the vector  $V = V(D)$ ; if there is a  $\alpha$ -significant element in the implicit vector  $V = V(D)$  then  $\text{Certificate}(D) = 1$ .

A function is a  $(\beta, \tau)$ -mimic if an existence of a  $\beta$ -significant element  $v_i$  in  $V = V(D)$  implies that  $\text{Mimic}(D) \in [(1 - \tau)|v_i|, (1 + \tau)|v_i|]$ .

Consider the case where we are given two such functions. Then there is a simple yet effective method for finding heavy hitters. We call this the **Purifier** method.

**ALGORITHM 2.1. PURIFIER**

1.  $c = \text{Certificate}_{V, \alpha, \beta}(D)$ .
2.  $\mu = \text{Mimic}_{V, \beta}(D)$ .
3. If  $c = 1$  output  $\mu$ ; else output 0.

It is easy to see that **Purifier** solves the following promise problem. If a vector has an  $\alpha$ -significant element then the output will be its approximation. If, on the other hand, there is no  $\beta$ -significant element, the output is 0. Generally speaking, we are able to find extremely heavy elements and to discard somewhat smaller ones.

It is a long way from this idea to a specific solution. To apply **Purifier**, we need to “guess” the right parameters, the certificate and the mimic function. In addition, we need to prove the correctness of the certificate and the mimic function, and prove memory bound for their computations. It is also important to note that mimic and certificate are problem specific. There might be other technical issues such as dealing with the probabilistic nature of the algorithms. The above tasks may be highly non-trivial. Still, the idea of **Purifier** provides a general framework for streaming algorithms; thus we believe it will be helpful for future applications.

We stress that Indyk and Woodruff were the first to implicitly use the **Purifier** method to tackle  $F_k$  heavy hitters. They observed that  $F_2$  can be used as a certifier for  $F_k$  heavy hitters. In this paper we generalize this idea and show that  $F_2$  is a litmus task of tractability; that is, any tractable frequency-based function can be both certified and mimicked by  $F_2$ . (From this, alternative perspective, in [13] we show that  $\log(n)$ -approximations can certify  $\epsilon$ -heavy rows, and we give another, completely different function for mimicking.)

**Purifier** can be seen as a reduction from our target function to a mimic function under specific conditions. It gives a hierarchy of functions that can be reduced. We remark that it will also be very interesting to construct such complexity hierarchies for more general settings.

<sup>5</sup>It is important to note that heavy hitters were solved for many specific metrics [18, 22].

### 3. THE LOWER BOUND

To establish lower bounds, we will use SET DISJOINTNESS and INDEX problems from communication complexity [37]. Recall that SET DISJOINTNESS is the following promise problem: each of  $t \geq 2$  players is given a set from the universe  $[N]$ ; all sets have exactly one common element or disjoint. The lower bound on the communication complexity of this problem is  $\Omega(\frac{N}{t})$  for the randomized one-way communication complexity [17]. INDEX is the following promise problem: there are two players: Alice and Bob. Alice is given a set  $S$  from the universe  $[N]$ ; Bob is given a single element  $x \in [N]$ . The players must decide whether or not  $x \in S$ . The lower bound on the communication complexity of INDEX problem is  $\Omega(N)$  for the randomized one-way communication complexity [37]. We stress that SET DISJOINTNESS and INDEX are common tools to prove lower bounds on streams. We start with the following observation.

**FACT 3.1.** *Let  $G$  be a non-decreasing function such that  $G(1) = 0$ . Then  $G \notin \text{STREAM-POLYLOG}$ .*

**PROOF.** Consider any algorithm that solves  $G\text{-Sum}(D, 2)$ . Let  $y = \min\{z \in \mathbb{N}^+ : G(z) > 0\}$ . We have  $y > x_0 \geq 1$ . Consider the following reduction to the SET DISJOINTNESS problem with two players. Alice repeats each element of his set  $y - 1$  times. Then she evaluates the algorithm for  $G\text{-Sum}(D, 2)$  on the stream of repetitions  $D$  and sends the resulting memory to Bob. Bob continues its computations by feeding the same algorithm his input. The algorithm may have two outputs: 0 or 2-approximation of  $G(y) > 0$ . Clearly, it must distinguish between these two cases. Thus, any algorithm that solves  $G\text{-Sum}(D, 2)$  must use a linear memory.  $\square$

In the remainder of this paper we concentrate on the case where  $G(x) > 0$  for all positive integers  $x$ . Further, we may assume that

$$G(1) = 1. \quad (2)$$

Indeed otherwise, we shall consider  $G'(x) = \frac{G(x)}{G(1)}$ , approximate  $G'(M(D))$  and multiply by  $G(1)$ ; the result is an approximation of  $G(M)$ .

**THEOREM 3.2.** *The Lower Bound for Non-Tractable Functions.*

*Let  $G$  be a non-decreasing function such that  $G(0) = 0$ . If  $G$  is not tractable, then  $G \notin \text{STREAM-POLYLOG}$ .*

**PROOF.** By Fact 3.1, it is sufficient to consider the case when  $G(1) = 1$  and (1) does not hold. In particular:

$$\begin{aligned} &\exists k \forall N_0 \forall t \exists x, y \in \mathbb{N}^+, \exists R \in \mathbb{R}^+ \exists \epsilon : \\ &\left( R > N_0, \frac{G(x)}{G(y)} = R, \epsilon > \frac{1}{\log^k(Rx)} \right) \cap \\ &\left( \left( \frac{\pi_\epsilon(x)}{y} \right)^2 < \frac{R}{\log^t(xR)} \right). \end{aligned} \quad (3)$$

Consider the fixed  $k$  from (3) and let  $t, N_0$  be any arbitrary large numbers. There exists  $R \geq N_0$ ; a pair  $(x, y)$  and  $\epsilon > \log^{-k}(Rx)$  such that (3) holds. Denote  $N = \lfloor R \rfloor$ ; then  $N \leq \frac{G(x)}{G(y)} < N + 1$  and  $\left( \frac{\pi_\epsilon(x)}{y} \right)^2 < \frac{N+1}{\log^t(xN)}$ .

First, consider the case when  $\pi_\epsilon(x) = x > y$ ; in this case we show that even a small constant approximation is not

possible in polylogarithmic space. Consider a SET DISJOINTNESS on the domain of size  $\Theta(N)$  with  $s = \lceil \frac{x}{y} \rceil$  players. Consider the following reduction: the first player repeats  $y$  times every element of his subset, applies the streaming algorithm to the resulting stream, and sends the memory content to the second player. The second player and the rest of the players repeat their sets  $y$  times and continue the execution of the algorithm on the resulting stream. Consider the resulting stream  $D(n, m)$ , where  $n = N$  and  $m \leq n(x+1)$ . If the sets are disjoint, then the resulting value of  $G$ -Sum should be  $NG(y)$ . If there is an intersection, then the resulting value should be at least  $((N-1)G(y) + G(x)) \approx 2NG(y)$  for sufficiently large  $N$ . Thus, there exists a small constant  $c$  such that a  $c$ -approximation of  $G$ -Sum solves the SET DISJOINTNESS problem. Thus, any algorithm for  $G$ -Sum( $D, c$ ) must use  $\Omega(\frac{N}{s^2})$  memory bits. However, (3) implies:

$$\frac{N}{s^2} \geq 0.5 \log^t(xN) \geq 0.5(0.5 \log(nm))^t. \quad (4)$$

Next, consider the case that  $\pi_\epsilon(x) < x$ . W.l.o.g., assume that  $G(x + \pi_\epsilon(x)) > (1 + \epsilon)G(x)$ . Assume that  $\pi_\epsilon(x) \leq y$  and consider the following reduction to the INDEX problem of dimensionality  $N$ : Alice repeats elements of her set  $y$  times and Bob repeats his index  $x$  times. There are two possible outcomes for  $G(M(D))$ : either  $G(M(D)) = NG(y) + G(x)$  or  $G(M(D)) = (N-1)G(y) + G(x + \pi_\epsilon(x)) > (N-1)G(y) + (1 + \epsilon)G(x)$ . Since  $G(x) \geq NG(y)$  it follows that there exists  $\epsilon' = \Omega(\epsilon)$  such that any algorithm that solves  $G$ -Sum( $D, \epsilon'$ ) must solve INDEX. Thus, the lower bound is  $N \geq \frac{N}{s^2}$ ; i.e., we obtain (4).

Finally, consider the case that  $y < \pi_\epsilon(x) < x$ . Consider a SET DISJOINTNESS problem with  $s = \lceil \frac{\pi_\epsilon(x)}{y} \rceil$  players over a domain of size  $N$ . Consider the same reduction as in the case  $x = \pi_\epsilon(x)$  with an additional modification. The last player does the following: for each element  $a$  of his set he simulates separately the insertion to the stream  $x$  copies of  $a$ . Thus, he simulates execution of the algorithm on streams  $D_a = D', a, \dots, a$ , where  $D'$  is the stream generated by the first  $s-1$  players and  $a$  is an element from the  $s$ -th player set. Consider the outputs of  $G$ -Sum on all these streams. If the sets are disjoint then the output will be  $N'G(y) + G(x)$  for all streams for some fixed  $N' \leq N$ . If there is an intersection, then exactly one answer will be  $N'G(y) + G(x + \pi_\epsilon(x)) > N'G(y) + (1 + \epsilon)G(x)$ . Since  $G(x) \geq NG(y)$ , we conclude that there exists  $\epsilon' = \Omega(\epsilon)$  such that any algorithm that solves  $G$ -Sum( $D, \epsilon'$ ) will distinguish between these two results. By repeating the above arguments, we also obtain (4). Finally, the case when  $G(x - \pi_\epsilon(x)) < (1 - \epsilon)G(x)$  can be handled similarly; instead of working with  $x$  we use  $x - \pi_\epsilon(x)$ .

To summarize, we have shown that there exists  $k$  such that for any  $t$  there exists  $n, m, \epsilon \geq \log^k(nm)$  and a stream  $D = D(n, m)$  such that any algorithm that solves  $G$ -Sum( $D, \epsilon$ ) requires  $\Omega(\log^t(nm))$  memory bits. Thus, and by Definition 1.2 we proved that  $G \notin \text{STREAM-POLYLOG}$ .  $\square$

## 4. DEFINITIONS, NOTATIONS AND COMMENTS

In this section we summarize definitions and notations that we need to establish the upper bound. A probability  $p$  is *asymptotically equal* to  $a$  if  $|p - a| = O^*(1/(nm))$ ; we

denote this fact by  $p \approx a$ . A probability  $p$  is *negligible* if  $p \approx 0$ . We will extensively use the fact that a union of a polylogarithmic number of negligible events is a negligible event (i.e., occurs with negligible probability).

### DEFINITION 4.1. Major elements

Let  $V$  be a vector of dimensionality  $n$  with entries  $v_i$ . Let  $G$  be a function and  $d$  be a positive real number such that  $d \geq 1$ . An element  $v_i$  is a  $(G, d)$ -major with respect to  $V$  if:

$$G(v_i) > d \sum_{j \neq i} G(v_j).$$

### DEFINITION 4.2. Sampled stream

Let  $D$  be a stream and  $H : [n] \mapsto \{0, 1\}$  be a function. A sampled stream  $D_H$  is a stream defined as  $D \cap H^{-1}(1)$ ; i.e.,  $D_H$  contains elements of  $D$  that are mapped to 1 by  $H$ .

### DEFINITION 4.3. Residual second moment.

For a vector  $V$  with  $(F_2, 1)$ -major entry  $v_i$ , define  $F_2^{\text{res}}(V) = F_2(V) - v_i^2$ . I.e.,  $F_2^{\text{res}}(V)$  is a second moment of  $V$  minus the square of the maximal element.

### DEFINITION 4.4. $G$ -Vector

Let  $G$  be a function and  $V$  be a vector of dimensionality  $n$  with entries  $v_i, i \in [n]$ . A  $G$ -vector  $G(V)$  is a vector of dimensionality  $n$  with  $i$ -th entry equal to  $G(v_i)$ .

### DEFINITION 4.5. Hadamard product

For two vectors of dimensionality  $n$  define  $\text{Had}(V, U)$  to be their Hadamard product; i.e.,  $\text{Had}(V, U)$  is a vector of dimensionality  $n$  with entries  $v_i u_i$ .

Every vector  $H$  of dimensionality  $n$  with entries  $h_i \in R$  defines a function  $H : [n] \mapsto R$  as  $H(i) = h_i$ . We thus will freely interchange the notions of random vector and hash functions. In our auxiliary algorithms we use notations  $D, M(D), \epsilon$ , etc., to define an input of the auxiliary algorithm. It is important to distinguish between the parameters of auxiliary algorithms and the ones of the  $G$ -Sum problem.

## 5. COMPUTING $G$ -CORE

In this section we solve the following problem:

$G$ -Core( $D, \alpha$ )

Given: A stream  $D$ , and a real  $\alpha$ .

Output: A set of positive numbers  $S = \{s_1, \dots, s_l\}$  such that  $l = O^*(1)$  and:

1. There exist a sequence  $j_1 < \dots < j_l$  such that  $s_i$  is a  $(1 \pm \alpha)$ -approximation of a  $G(m_{j_i})$ .
2. If there is a  $(G, 1)$ -major element  $m_i$  w.r.t.  $M(D)$ , then  $i \in \{j_1, \dots, j_l\}$ , i.e.,  $S$  contains a  $(1 \pm \alpha)$ -approximation of a  $G(m_i)$ .

### 5.1 Preliminaries

FACT 5.1. Let  $V$  be a vector with non-negative entries. Let  $H$  be a pairwise-independent random vector (of the same dimensionality as  $V$ ) with zero-one entries  $h_i$  that are uniformly distributed. Let  $H'$  be a vector with entries  $1 - h_i$ .

Denote  $X = \langle V, H \rangle$  and  $Y = \langle V, H' \rangle$ , where  $\langle \cdot, \cdot \rangle$  is an inner product. If for all  $i$ ,  $v_i < 0.01|V|$  then  $P((X < \frac{1}{3}|V|) \cup (Y < \frac{1}{3}|V|)) < 0.1$ .

PROOF. Clearly,  $X = |V| - Y$  and  $E(X) = E(Y) = 0.5|V|$ . Further, since  $H$  is pairwise independent and by the condition of the lemma, we have that  $\text{Var}(X) = 0.25F_2(V) \leq \frac{1}{400}|V|^2$ . Thus, by Chebyshev inequality:  $P((X < \frac{1}{3}|V|) \cup (Y < \frac{1}{3}|V|)) = P(|X - E(X)| \geq \frac{1}{6}|V|) \leq \frac{36}{|V|^2} \frac{|V|^2}{400} < 0.1$ .  $\square$

LEMMA 5.2. Let  $V$  be a vector with non-negative entries. Let  $H$  be a pairwise independent random vector (of the same dimensionality as  $V$ ) with zero-one entries  $h_i$  that are uniformly distributed. Let  $H'$  be a vector with entries  $1 - h_i$ . Denote  $X = \langle V, H \rangle$  and  $Y = \langle V, H' \rangle$ . Let  $K$  be a parameter,  $K > 10^4$ . If there exists  $(L_1, K)$ -major element  $v_i$  then:

$$P((X > KY) \vee (Y > KX)) = 1.$$

If there is no  $(L_1, 10^{-4}K)$ -major element then:

$$P((X > KY) \vee (Y > KX)) \leq 0.5.$$

PROOF. The first claim of the lemma follows directly. Indeed, w.l.o.g., assume that  $v_1 > K \sum_{i>1} v_i$  and that  $H(1) = 1$ . Then  $X \geq v_1 > K \sum_{i>1} v_i \geq KY$ .

To show the second claim, consider an entry  $v_i$  with maximal value. If  $v_i < 0.01|V|$ , we are done by Fact 5.1; thus we assume that  $v_i \geq 0.01|V|$ . Assume that there exists  $v_j, j \neq i$  such that  $v_j \geq \frac{100}{K}v_i$ . We have:

$$Kv_j \geq 100v_i \geq |V|, \quad Kv_i \geq 0.01K|V| > |V|. \quad (5)$$

W.p. 0.5,  $h_i \neq h_j$ , in which case by (5) the event  $(X > KY) \vee (Y > KX)$  cannot happen.

Finally, assume that for all  $j \neq i$ , we have:

$$v_j < \frac{100}{K}v_i < 0.01(|V| - v_i). \quad (6)$$

Consider vector  $V'$  that is equal to  $V$  for all entries and is equal to 0 on its  $i$ -th entry. Consider  $X' = \langle V', H \rangle$  and  $Y' = \langle V', H' \rangle$ . By (6) and by Fact 5.1:  $P((X' < \frac{1}{3}|V'|) \vee (Y' < \frac{1}{3}|V'|)) \leq 0.1$ . But this implies that w.p. at least 0.9:  $KX \geq KX' \geq \frac{K}{3}(|V| - v_i) \geq |V|$ ,  $KY \geq KY' \geq \frac{K}{3}(|V| - v_i) \geq |V|$ , in which case the event  $(X > KY) \vee (Y > KX)$  cannot happen.  $\square$

LEMMA 5.3. Let  $V$  be a vector with non-negative entries. Let  $H$  be a random vector (of the same dimensionality as  $V$ ) with pairwise independent entries  $h_i \sim U(\{-1, 1\})$ . Denote  $X = \langle V, H \rangle$ . If there exists  $i$  such that  $v_i$  is  $(F_2, 1)$ -major with respect to  $V$  then:

$$P(|X| - v_i \geq 2(F_2^{res}(V))^{0.5}) \leq 0.25. \quad (7)$$

PROOF. Denote  $Z = \sum_{j \neq i} h_j v_j$ . By triangle inequality:  $|X| \leq |Z| + |v_i|$ . By inverse triangle inequality:  $|X| = |v_i - (-Z)| \geq |v_i| - |Z|$ . Since  $v_i$  is non-negative, we have:  $v_i - |Z| \leq |X| \leq v_i + |Z|$ .

By the linearity of expectation,  $E(Z) = 0$ . By pairwise independence of  $H$ , it follows that  $\text{Var}(Z) = F_2(V) - v_i^2 = F_2^{res}$ . Thus, by Chebyshev inequality:  $P(|Z| \geq 2\sqrt{F_2^{res}(V)}) \leq 0.25$ .

$\square$

The following fact is folklore; a similar statement can be found, e.g., in [7].

FACT 5.4. Let  $D$  be a stream such that  $M(D)$  contains a  $(F_2, 2)$ -major element  $m_i$ . There exists an algorithm that makes a single pass over  $D$ , uses a polylogarithmic memory and outputs  $r$  such that:  $2(F_2^{res}(M))^{0.5} < r < 3(F_2^{res}(M))^{0.5}$ . The algorithm errs with a negligible probability.

PROOF. W.l.o.g, assume that  $m_1$  is  $(F_2, 2)$ -major element. Let  $H$  be a zero-one hash function with pairwise independent entries. Consider two streams  $D_H$  and  $D_{1-H}$  and apply the AMS algorithm for  $F_2$  on both streams. We require that the outputs be 1.1-approximations of the second moments with a negligible probability of error. Let  $a$  and  $b$  be two outputs. Thus  $Z = 10\min\{a, b\}$  is a 1.1 approximation of the following random variable:

$$X = 10 \sum_{i>1} \mathbf{1}_{H(i) \neq H(1)} v_i^2.$$

Since  $H$  is pairwise independent we have that  $E(X) = 5F_2^{res}(M)$ . Since  $0 \leq X \leq 10F_2^{res}(V)$ , it follows that  $\text{Var}(X) \leq E(X^2) \leq (10F_2^{res}(V))^2$ . Let  $Y$  be an average of  $C = O(1)$  independent  $X$ ; it follows that  $\text{Var}(Y) \leq C^{-1}(10F_2^{res}(V))^2$ . Thus for sufficiently large  $C$  by Chebyshev inequality:

$$P(|Y - 5F_2^{res}| \geq 0.1F_2^{res}) \leq 0.1. \quad (8)$$

Taking a median of  $O(\log(nm))$  independent averages drops the probability to negligible. Let  $Q$  be a median of  $O(\log(nm))$  averages of  $C$  independent  $Z$ s. By the union bound all  $Z$  will be 1.1 approximations of corresponding  $X$ s except with negligible probability. Thus and by (8),  $r = Q^{0.5}$  satisfies the conditions of the fact except with negligible probability.  $\square$

FACT 5.5. Let  $\epsilon \leq 0.5$  and let  $x, a, b \geq 0$  be such that  $|x - a| \leq 0.1\pi_\epsilon(x)$  and  $b < 0.1\pi_\epsilon(x)$ . Then  $(1 - 4\epsilon)G(a + b) \leq G(a) \leq (1 + 4\epsilon)G(a - b)$ .

PROOF. We have  $a - b \geq x - 0.2\pi_\epsilon(x) > 0$  and  $a + b \leq x + 0.2\pi_\epsilon(x)$ . Thus  $(1 - \epsilon)G(x) \leq G(a - b) \leq G(a) \leq G(a + b) \leq (1 + \epsilon)G(x)$ . Thus  $(1 + 4\epsilon)G(a - b) \geq \frac{(1 - \epsilon)(1 + 4\epsilon)}{(1 + \epsilon)}G(a) \geq G(a)$ . The second part can be proven similarly.  $\square$

FACT 5.6. Let  $x, a, b$  be such that  $|x - a| \leq b$ . Then  $(1 - \epsilon)G(a + b) \leq G(a) \leq (1 + \epsilon)G(a - b) \implies (1 - \epsilon)G(x) \leq G(a) \leq (1 + \epsilon)G(x)$ .

PROOF. Assume, on the contrary, that  $G(a) > (1 + \epsilon)G(x)$ . Then since  $x \geq a - b$  we have a contradiction:  $G(a) > (1 + \epsilon)G(x) \geq (1 + \epsilon)G(a - b) \geq G(a)$ . The second condition is proven similarly.  $\square$

## 5.2 Algorithms for $G$ -Core

We define and solve the following problem:

Hybrid-Major( $D, \epsilon$ )

Given: A stream  $D$ , and a positive real  $\epsilon$ .

Output: A real  $r \geq 0$  such that:

1. If  $r \neq 0$  then  $r$  is an  $(1 \pm 4\epsilon)$ -approximation of a  $G(m_j)$  for some  $m_j$ .
2. If there exists a  $(F_2, 1)$ -major element  $m_i$  (w.r.t.  $M(D)$ ) such that  $\pi_\epsilon(m_i) \geq 20^5(F_2^{res}(M(D)))^{0.5}$  then  $r$  is a  $(1 \pm 4\epsilon)$ -approximation of  $G(m_i)$ .

<sup>6</sup>In general, the inequality may not hold if  $\pi_\epsilon(x) = 1$  since Definition 1.3 limits  $\pi$  to integers. However, in this case either  $G$  is not tractable or  $x = O^*(1)$  and can be approximated precisely. We thus assume that  $\pi_\epsilon(x) \neq 1$ .

ALGORITHM 5.7. *Compute-Hybrid-Major*( $D, \epsilon$ )

1. Repeat  $O(\log(nm))$  times, independently and in parallel:
  - (a) Generate uniform pairwise independent vector  $H' \in \{0, 1\}^n$  and compute  $a_t = |\langle M, H' \rangle|$ .
2. Compute  $a = \text{median}\{a_t\}$ .
3. Repeat  $O(\log(nm))$  times, independently and in parallel:
  - (a) Generate uniform pairwise independent vector  $H \in \{0, 1\}^n$ .
  - (b) Using the AMS algorithm<sup>a</sup>, compute 5/4-approximations of second moments:  
 $X' = F_2(D_H)$  and  $Y' = F_2(D_{(1-H)})$  with negligible probability of error.
  - (c) If  $X' < (20)^4 Y'$  and  $Y' < (20)^4 X'$  then output 0 and terminate the algorithm.
4. In parallel, apply the algorithm for the residual moment from Fact 5.4. Let  $b$  be the output of the algorithm.
5. If  $(1-4\epsilon)G(a+b) > G(a)$  or  $G(a) > (1+4\epsilon)G(a-b)$  then output 0.
6. Otherwise output  $G(a)$ .

<sup>a</sup>The sketching algorithm for  $F_2$  approximation from [2].

LEMMA 5.8. *Algorithm Compute-Hybrid-Major solves Hybrid-Major( $D, \epsilon$ ) with negligible probability of error.*

PROOF. First, we show that if there is no  $(F_2, 2)$ -major entry, then the output is 0 except with negligible probability. Consider a single iteration of the main loop of the algorithm. Consider a vector  $M'$  with entries  $m_i^2$  and denote  $X = \langle M', H \rangle, Y = |M'| - \langle M', H \rangle$ . By the properties of the AMS algorithm, with negligible probability of error,  $X'$  is 5/4-approximation of  $X$  and  $Y'$  is 5/4-approximation of  $Y$ . By Lemma 5.2, this implies that w.p. at most to  $0.5 + o(1)$ :

$$X' \leq \frac{5}{4}X \leq \frac{5}{2}(10)^4 Y < (20)^4 Y', \quad Y' < (20)^4 X'.$$

Thus, except with negligible probability, the algorithm will output 0.

Assume that there is a  $(F_2, 2)$ -major entry  $m_i$ . Then by Lemma 5.3 and by Chernoff bound,  $|m_i - a| \leq 2(F_2(M(D)))^{0.5}$  except with negligible probability. From Fact 5.4 it follows that  $2(F_2^{res}(M(D)))^{0.5} < b < 3(F_2^{res}(M(D)))^{0.5}$  except with negligible probability. Thus, by Fact 5.6 it follows that if the algorithm outputs  $G(a)$  then  $G(a)$  is  $(1 \pm 4\epsilon)$ -approximation of  $G(m_i)$ . Thus, the first condition of Hybrid-Major follows.

Finally, assume that  $\pi_\epsilon(m_i) \geq (20)^5 (F_2^{res}(M(D)))^{0.5}$ . Definition 1.3 implies that  $m_i \geq \pi_\epsilon(m_i)$  and thus  $m_i$  is  $(F_2, 10^5)$ -major w.r.t.  $M$ . Thus, by Lemma 5.2, we have (except with negligible probability):  $X' > 20^4 Y'$  or  $Y' > 20^4 X'$ . Thus, except with negligible probability, the algorithm will not terminate before the last line. Also,  $|m_i - a| \leq 2F_2^{res}(M) <$

$0.01\pi_\epsilon(m_i)$ , and  $b \leq 3F_2^{res}(M) < 0.01\pi_\epsilon(m_i)$ . Thus, by Fact 5.5 the algorithm will output  $G(a)$  which is a  $(1 \pm 4\epsilon)$ -approximation of  $G(m_i)$ . Thus, the second condition of Hybrid-Major follows.  $\square$

Lemma 5.10 states that for any  $(G, 1)$ -major entry of the vector  $M(D)$ , the square of its local jump  $\pi_\epsilon$  is “heavy” with respect to  $F_2^{res}(M)$  (at least after deleting a small number of entries). First, we need to establish the following simple corollary of tractability.

FACT 5.9. *If  $G$  is tractable then*

$$\exists N_1 \forall N > N_1 \in \mathbb{N}^+ : G(N) \leq N^3. \quad (9)$$

PROOF. Indeed, if (9) does not hold then  $G$  cannot be tractable. In particular, there is an arbitrary large  $N$  such that  $G(N) > N^3$  and thus  $N^2 < \log^{-t}(NG(N))G(N)$  for any  $t = O(1)$ . By Definition 1.4,  $G$  is not tractable (for  $k = 0$ , arbitrary large  $t$  and  $x = N, y = 1$ ).  $\square$

LEMMA 5.10. *Let  $G$  be a non-decreasing tractable function. Then for any  $k = O(1)$  there exists  $t = O(1)$  such that for any  $n, m$  and for any  $\epsilon > \log^{-k}(nm)$  the following is true. Let  $D = D(n, m)$  be a stream that defines frequency vector  $M = M(D)$ . If there exists a  $(G, 1)$ -major element  $m_i$  w.r.t.  $M$ , then there exists set  $S \subseteq [n]$  such that  $|S| = O(\log(m))$  and:*

$$(\pi_\epsilon(m_i))^2 = \Omega(\log^{-t}(nm) \sum_{j \notin S \cup \{i\}} m_j^2).$$

PROOF. W.l.o.g., assume that  $m_1$  is the  $(G, 1)$ -major entry. Consider layers  $S_s = \{j \in [n] : 2^s \leq m_j < 2^{s+1}\}$  for  $s = 0, 1, \dots$ . Let  $N_0$  be the constant from Definition 1.4; and  $N_1$  be the constant from Fact 9 (i.e.,  $G(x) \leq x^3$  for  $x > N_1$ ). Let  $W = \{s : |S_s| \leq N_0 + N_1\}$  and let  $S = \cup_{s \in W} S_s$ . Note that  $m_i \leq m$  for all  $i \in [n]$ ; thus there are at most  $\lceil \log(m) \rceil$  non-empty layers. Thus  $|S| = O(\log(m))$ , i.e.,  $S$  satisfies the condition of the lemma.

Consider  $X = \sum_{j \neq i, j \notin S} m_j^2$ . If  $X = 0$  then the lemma follows. Otherwise there exists at least one layer  $S_l, l \notin W$  such that  $\sum_{j \in S_l} m_j^2 \geq \frac{1}{\lceil \log(m) \rceil} X$ . Since  $m_1$  is  $(G, 1)$ -major w.r.t.  $M$  and by monotonicity of  $G$ , it follows that  $G(m_1) \geq \sum_{j > 1} G(m_j) \geq |S_l| G(2^l)$ . Since  $l \notin W$  we have that  $|S_l| > N_0$ ; thus  $\frac{G(m_1)}{G(2^l)} \geq |S_l| > N_0$ . Since  $G$  is tractable, by Definition 1.4 of tractable functions there exists  $t = t(k) = O(1)$  such that:  $\left(\frac{\pi_\epsilon(m_i)}{2^l}\right)^2 \geq \frac{1}{\log^t\left(\frac{G(m_1)}{G(2^l)}\right)} \left(\frac{G(m_1)}{G(2^l)}\right)$ . Also, since

$G(m_1) \geq |S_l| > N_1$  we conclude that  $G(m_1) \leq m_1^3 \leq m^3$ ; thus  $\log^t\left(\frac{G(m_1)}{G(2^l)}\right) \leq (3 \log(m))^t$ . Thus, we summarize:  $\left(\frac{\pi_\epsilon(m_i)}{2^l}\right)^2 \geq \frac{1}{(3 \log(m))^t} |S_l|$ ; and finally:

$$(\pi_\epsilon(m_i))^2 \geq \frac{1}{(3 \log(m))^t} 2^{2l} |S_l| \geq \frac{1}{(12 \log(nm))^{t+1}} X.$$

$\square$

Consider the following algorithm that solves  $G$ -Core.

ALGORITHM 5.11. *Compute-G-Core*( $D, \epsilon, p$ )

1. Generate a pairwise independent hash function  $H : [n] \mapsto \tau$ ,  
where  $\tau = O^*\left(\frac{1}{p\epsilon}\right)$ .
2. For each  $k \in [\tau]$  compute in parallel  $c_i = \text{Compute-Hybrid-Major}(D_{H_k}, 0.25\epsilon)$ ,  
where  $H_k(i) = \mathbf{1}_{H(i)=k}$ .
3. Output  $S = \{G(c_i) : c_i > 0\}$ .

Finally, let us state the following claims which are a generalization of the arguments from [33] and can be seen as a variant of Lemmas 5.3 and 5.5 from our other paper [13]. (It is important to note, however, that while these results have similarities with the lemmas from [13], they are different. In this paper we consider frequency vectors, where in [13] we work with implicit tensors). We refer the reader to the full version [10] for details.

**THEOREM 5.12.** *Algorithm Compute-G-Core solves  $G\text{-Core}(D, \epsilon)$  and errs w.p. asymptotically equal to  $p$ . Compute-G-Core uses  $O^*(1)$  memory bits if  $p = \Omega^*(1)$  and  $\epsilon = \Omega^*(1)$ .*

**PROOF.** First, except with negligible probability, every positive  $c_i$  is  $(1 \pm \epsilon)$ -approximation of some distinct entry  $G(m_j)$ . Second, assume that there exists a  $(G, 1)$ -dominant entry  $m_i$ . Denote  $X = \sum_{j \neq i} v_i^2 \mathbf{1}_{H(j)=H(i)}$ . By pairwise independence of  $H$ , we have  $E(X) = \frac{1}{\tau}(F_2(M) - m_i^2)$ . By Lemma 5.10 there exists a set  $S$  such that  $|S| = O^*(1)$  and:

$$\pi_\epsilon(m_i)^2 = \Omega^*\left(\sum_{j \notin S} m_j^2\right). \quad (10)$$

Let  $\mathcal{A}$  be an event that  $\pi_\epsilon(m_i)^2 > (20)^5 X$  and  $H(j) \neq H(i)$  for any  $j \in S$ . By Markov inequality, by pairwise independence of  $H$  and by (10), there exists  $\tau = \Omega^*\left(\frac{1}{p}\right)$  such that:  $P(\mathcal{A}) \geq (1 - p)$ . If  $\mathcal{A}$  occurs, then by Lemma 5.8  $c_{H(i)}$  is  $(1 \pm \epsilon)$ -approximation of  $G(m_i)$  except with negligible probability. Thus, the final probability of error is approximately equal to  $p$ .

Algorithm *Compute-G-Core* can be seen as a computational tree of a constant depth (after substituting the auxiliary algorithms). Each internal node in the tree has a polylogarithmic number of children. Each leaf is either a direct computation on a stream that requires polylog space or AMS algorithm for  $F_2$  that also requires polylog space. Thus, the total space is polylogarithmic. In fact, for any  $k = O(1)$  there exists  $t = t(k) = O(1)$  such that we can solve  $G\text{-Core}(D(n, m), \epsilon)$ , where  $\epsilon \geq \log^{-k}(nm)$  with  $O(\log^{-t}(nm))$  space.  $\square$

**THEOREM 5.13.** *If there exists an algorithm that solves  $G\text{-Core}$  using memory  $O^*(1)$  and a single pass over  $D$  with probability of error  $\Omega^*(1)$ , then there exists an algorithm that solves  $G\text{-sum}$  using memory  $O^*(1)$  and a single pass over  $D$  with probability of error bounded at most  $1/3$ .*

### 5.3 Proof of the Main Theorem 1.5

**PROOF.** Theorem 3.2 proves that tractability is necessary. The sufficient condition, i.e., the existence of an algorithm that solves  $G\text{-Sum}$  for tractable  $G$ , follows from Theorem 5.13 and Theorem 5.12.  $\square$

## 6. ACKNOWLEDGEMENT

We thank Piotr Indyk for a helpful discussion and for pointing out Question 5 in [40].

## 7. REFERENCES

- [1] C. C. Aggarwal. *Data Streams: Models and Algorithms (Advances in Database Systems)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [2] N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.*, 58(1):137–147, 1999.
- [3] Z. Bar-Yossef, T. S. Jayram, R. Kumar, and D. Sivakumar. An information statistics approach to data stream and communication complexity. *J. Comput. Syst. Sci.*, 68(4):702–732, 2004.
- [4] Z. Bar-Yossef, T. S. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream. In *RANDOM '02: Proceedings of the 6th International Workshop on Randomization and Approximation Techniques*, pages 1–10, London, UK, 2002. Springer-Verlag.
- [5] P. Beame, T. S. Jayram, and A. Rudra. Lower bounds for randomized read/write stream algorithms. In *STOC '07: Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 689–698, New York, NY, USA, 2007. ACM.
- [6] L. Bhuvanagiri and S. Ganguly. Estimating entropy over data streams. In *ESA '06: Proceedings of the 14th conference on Annual European Symposium*, pages 148–159, London, UK, 2006. Springer-Verlag.
- [7] L. Bhuvanagiri, S. Ganguly, D. Kesh, and C. Saha. Simpler algorithm for estimating frequency moments of data streams. In *SODA*, pages 708–713, 2006.
- [8] V. Braverman, K.-M. Chung, Z. Liu, M. Mitzenmacher, R. Ostrovsky, “AMS Without 4-Wise Independence on Product Domains,” STACS 2010.
- [9] V. Braverman, R. Ostrovsky, “Effective Computations on Sliding Windows,” SIAM J. Comput. Volume 39, Issue 6, pp. 2113-2131 (2010).
- [10] V. Braverman, R. Ostrovsky, “Zero-One Frequency Laws,” ArXiv, 2010.
- [11] V. Braverman, R. Ostrovsky, “Smooth Histograms for Sliding Windows,” In Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science (October 21 - 23, 2007).
- [12] V. Braverman, R. Ostrovsky, C. Zaniolo, “Optimal sampling from sliding windows,” In Proceedings of the Twenty-Eighth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (Providence, Rhode Island, USA, June 29 - July 01, 2009).
- [13] V. Braverman, R. Ostrovsky, “Measuring Independence of Datasets,” STOC 2010.
- [14] A. Chakrabarti, K. D. Ba, and S. Muthukrishnan. Estimating entropy and entropy norm on data streams. In *In Proceedings of the 23rd International Symposium on Theoretical Aspects of Computer Science STACS 2006*. Springer, 2006.
- [15] A. Chakrabarti, G. Cormode, and A. McGregor. A near-optimal algorithm for computing the entropy of a stream. In *SODA '07: Proceedings of the eighteenth*

- annual ACM-SIAM symposium on Discrete algorithms, pages 328–335, Philadelphia, PA, USA, 2007. Society for Industrial and Applied Mathematics.
- [16] A. Chakrabarti, G. Cormode, and A. McGregor. Robust lower bounds for communication and stream computation. In *STOC '08: Proceedings of the 40th annual ACM symposium on Theory of computing*, pages 641–650, New York, NY, USA, 2008. ACM.
- [17] A. Chakrabarti, S. Khot, and X. Sun. Near-optimal lower bounds on the multi-party communication complexity of set disjointness. In *IEEE Conference on Computational Complexity*, pages 107–117, 2003.
- [18] M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. In *ICALP '02: Proceedings of the 29th International Colloquium on Automata, Languages and Programming*, pages 693–703, London, UK, 2002. Springer-Verlag.
- [19] D. Coppersmith and R. Kumar. An improved data stream algorithm for frequency moments. In *SODA*, pages 151–156, 2004.
- [20] G. Cormode, M. Datar, P. Indyk, and S. Muthukrishnan. Comparing data streams using hamming norms (how to zero in). *IEEE Trans. on Knowl. and Data Eng.*, 15(3):529–540, 2003.
- [21] G. Cormode and M. Hadjieleftheriou. Finding frequent items in data streams. *Proc. VLDB Endow.*, 1(2):1530–1541, 2008.
- [22] G. Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *J. Algorithms*, 55(1):58–75, 2005.
- [23] J. Feigenbaum, S. Kannan, M. Strauss, and M. Viswanathan. An approximate  $\ell_2$ -difference algorithm for massive data streams. In *FOCS '99: Proceedings of the 40th Annual Symposium on Foundations of Computer Science*, page 501, Washington, DC, USA, 1999. IEEE Computer Society.
- [24] P. Flajolet and G. N. Martin. Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.*, 31(2):182–209, 1985.
- [25] S. Ganguly. Estimating frequency moments of data streams using random linear combinations. In *APPROX-RANDOM*, pages 369–380, 2004.
- [26] S. Ganguly and G. Cormode. On estimating frequency moments of data streams. In *APPROX '07/RANDOM '07: Proceedings of the 10th International Workshop on Approximation and the 11th International Workshop on Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 479–493, Berlin, Heidelberg, 2007. Springer-Verlag.
- [27] S. Guha, P. Indyk, and A. McGregor. Sketching information divergences. *Mach. Learn.*, 72(1-2):5–19, 2008.
- [28] S. Guha, A. McGregor, and S. Venkatasubramanian. Streaming and sublinear approximation of entropy and information distances. In *SODA '06: Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, pages 733–742, New York, NY, USA, 2006. ACM.
- [29] N. J. A. Harvey, J. Nelson, and K. Onak. Sketching and streaming entropy via approximation theory. In *FOCS '08: Proceedings of the 2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 489–498, Washington, DC, USA, 2008. IEEE Computer Society.
- [30] P. Indyk. Stable distributions, pseudorandom generators, embeddings, and data stream computation. *J. ACM*, 53(3):307–323, 2006.
- [31] P. Indyk and A. McGregor. Declaring independence via the sketching of sketches. In *SODA '08: Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 737–745, Philadelphia, PA, USA, 2008. Society for Industrial and Applied Mathematics.
- [32] P. Indyk and D. Woodruff. Tight lower bounds for the distinct elements problem. In *FOCS '03: Proceedings of the 44th Annual IEEE Symposium on Foundations of Computer Science*, page 283, Washington, DC, USA, 2003. IEEE Computer Society.
- [33] P. Indyk and D. Woodruff. Optimal approximations of the frequency moments of data streams. In *STOC '05: Proceedings of the thirty-seventh annual ACM symposium on Theory of computing*, pages 202–208, New York, NY, USA, 2005. ACM.
- [34] T. S. Jayram, A. McGregor, S. Muthukrishnan, and E. Vee. Estimating statistical aggregates on probabilistic data streams. In *PODS '07: Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 243–252, New York, NY, USA, 2007. ACM.
- [35] D. M. Kane, J. Nelson, and D. P. Woodruff. On the exact space complexity of sketching and streaming small norms. In *Proceedings of the 21st Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2010)*, 2010.
- [36] D. M. Kane, J. Nelson, and D. P. Woodruff. An Optimal Algorithm for the Distinct Elements Problem. *PODS 2010*, 2010.
- [37] E. Kushilevitz and N. Nisan. *Communication complexity*. Cambridge University Press, New York, NY, USA, 1997.
- [38] A. Lall, V. Sekar, M. Ogihara, J. Xu, and H. Zhang. Data streaming algorithms for estimating entropy of network traffic. *SIGMETRICS Perform. Eval. Rev.*, 34(1):145–156, 2006.
- [39] P. Li. Compressed counting. In *SODA '09: Proceedings of the Nineteenth Annual ACM -SIAM Symposium on Discrete Algorithms*, pages 412–421, Philadelphia, PA, USA, 2009. Society for Industrial and Applied Mathematics.
- [40] A. McGregor. Open problems in data streams and related topics. In *IITK workshop on algorithms for data streams*. <http://www.cse.iitk.ac.in/users/sganguly/data-stream-probs.pdf>, 2007., 2006.
- [41] S. Muthukrishnan. Data streams: algorithms and applications. *Found. Trends Theor. Comput. Sci.*, 1(2):117–236, 2005.
- [42] N. Nisan. Pseudorandom generators for space-bounded computations. In *STOC '90: Proceedings of the twenty-second annual ACM symposium on Theory of computing*, pages 204–212, 1990.
- [43] D. Woodruff. Optimal space lower bounds for all frequency moments. In *SODA '04: Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 167–175, 2004.