

Streaming Space Complexity of Nearly All Functions of One Variable on Frequency Vectors

Vladimir Braverman*
Johns Hopkins University
vova@cs.jhu.edu

Stephen R. Chestnut†
ETH Zurich
stephenc@ethz.ch

David P. Woodruff
IBM, Almaden Research Center
dpwoodru@us.ibm.com

Lin F. Yang†
Johns Hopkins University
lyang@pha.jhu.edu

ABSTRACT

A central problem in the theory of algorithms for data streams is to determine which functions on a stream can be approximated in sublinear, and especially sub-polynomial or polylogarithmic, space. Given a function g , we study the space complexity of approximating $\sum_{i=1}^n g(|f_i|)$, where $f \in \mathbb{Z}^n$ is the frequency vector of a turnstile stream. This is a generalization of the well-known frequency moments problem, and previous results apply only when g is monotonic or has a special functional form. Our contribution is to give a condition such that, except for a narrow class of functions g , there is a space-efficient approximation algorithm for the sum if and only if g satisfies the condition. The functions g that we are able to characterize include all convex, concave, monotonic, polynomial, and trigonometric functions, among many others, and is the first such characterization for non-monotonic functions. Thus, for nearly all functions of one variable, we answer the open question from the celebrated paper of Alon, Matias and Szegedy (1996).

Keywords

data streams, sketching, frequency moments

1. INTRODUCTION

One of the main open problems in the theory of data stream computation is to characterize functions of a frequency vector that can be computed, or approximated, efficiently. Here we characterize nearly all functions of the

*This material is based upon work supported in part by the National Science Foundation under Grant No. 1447639, by the Google Faculty Award and by DARPA grant N660001-1-2-4014. Its contents are solely the responsibility of the authors and do not represent the official view of DARPA or the Department of Defense.

†This material is based upon work supported in part by the National Science Foundation under Grant No. 1447639.

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PODS'16, June 26-July 01, 2016, San Francisco, CA, USA

© 2016 ACM. ISBN 978-1-4503-4191-2/16/06...\$15.00

DOI: <http://dx.doi.org/10.1145/2902251.2902282>

form $\sum_{i=1}^n g(|v_i|)$, where $v = (v_1, \dots, v_n) \in \mathbb{Z}^n$ is the frequency vector of the stream. This is a generalization of the famous frequency moments problem, where $g(x) = x^k$ for some $k \geq 0$, described by Alon, Matias, and Szegedy [1], who also asked:

“It would be interesting to determine or estimate the space complexity of the approximation of $\sum_{i=1}^n v_i^k$ for non-integral values of k for $k < 2$, or the space complexity of estimating other functions of the numbers v_i .” The first question was answered by Indyk and Woodruff [15] and Indyk [14] who were the first to determine, up to polylogarithmic factors, the space complexity of the frequency moments for all $k > 0$. We address the second question.

Braverman and Ostrovsky [6] and Braverman and Chestnut [5], answer the second question for sums of the form given above when g is monotone. We extend their characterizations to nearly all nonmonotone functions of one variable. Specifically, we characterize the set of functions g for which there exists a sub-polynomial $(1 \pm \epsilon)$ -approximation algorithm for the sum above. Our results can be adapted to characterize the set of functions approximable in polylogarithmic, rather than sub-polynomial, space. Among the main techniques we use is the layering method developed by Indyk and Woodruff [15] for approximating the frequency moments, and one may view our results as an exploration of the power of this technique.

Our results also partially answer an open question of Guha and Indyk at the IIT Kanpur workshop in 2006 [12] about characterizing sketchable distances - our results characterize nearly all distances that can be expressed in the form $d(u, v) = \sum_i g(|u_i - v_i|)$. Besides the aforementioned work on the frequency moments and monotone functions, other past work on this question was done by Guha, Indyk, and McGregor [13] as well as Andoni, Krauthgamer, and Razenshteyn [2]. Both of those papers address the same problem, but in domains different from ours. In the realm of general streaming algorithms, a result of Li, Nguyen, and Woodruff [20] shows that for any function g with a $(1 \pm \epsilon)$ -approximation algorithm implementable in sub-polynomial space, there exists a linear sketch that demonstrates its tractability. However, a drawback of that work is that the space complexity of maintaining the sketching matrix, together with computing the output, may be polynomially large. Our paper deals with a smaller class of functions than [20], but we provide a zero-one law and explicit algorithms.

While several of our lower bounds can be shown via re-

duction from the standard index and set-disjointness communication problems, to prove lower bounds for the widest class of functions we develop a new class of communication problems and prove lower bounds for them, which may be of independent interest. The problem, which we call **ShortLinearCombination**, is: given a finite set $\{0, a_1, \dots, a_r, b\}$ of possible frequencies, is there a frequency in the stream of value b ? Perhaps surprisingly, the communication complexity of **ShortLinearCombination** depends on the magnitudes of the coefficients in any linear combination expressing b in terms of a_1, \dots, a_r . We develop optimal communication bounds for this problem. This hints at the subtleties in proving lower bounds for a wide class of functions g , since by defining $g(b) \geq n \cdot \max_i g(a_i)$, the communication problem just described is a special case of characterizing the streaming complexity of all functions g . Our lower bounds provide a significant generalization of the set-disjointness problem, which has found many applications to data streams.

1.1 Applications

We now describe potential applications of this work.

1.1.1 Log-likelihood Approximation

One use for the approximation of non-monotonic functions is the computation of log-likelihoods. Here, the coordinates of the streamed vector v are taken to be i.i.d. samples from a discrete probability distribution. The likelihood of the vector v , under a distribution with probability mass function $p(x)$, $x \in \mathbb{Z}$, is $L(v) = \prod_{i=1}^n p(v_i)$, and its log-likelihood is $\ell(v) = -\sum_{i=1}^n \log p(v_i)$. Notice that $\ell(v)$ is of the form $\sum_{i=1}^n g(v_i)$, for $g(x) = -\log p(x)$. When $p(\cdot; \theta)$ is the distribution of a non-negative random variable or if $p(x; \theta)$ is symmetric about $x = 0$, our results can be applied to determine whether there is a streaming algorithm that efficiently approximates ℓ . In general, $-\log p(x)$ is not a monotonic function of x . For example, $p(x) = \lambda \frac{x^\alpha}{x!} e^{-\alpha} + (1 - \lambda) \frac{x^\beta}{x!} e^{-\beta}$ is a mixture of two Poissons, which is generally not monotonic. For any constants $\lambda, \alpha, \beta > 0$, the Poisson mixture log-likelihood $-\log p(x)$ satisfies our three criteria (to be described later), hence the log-likelihood of the stream can be approximated in poly-logarithmic space. Continuous distributions can be handled similarly by discretization.

Suppose that the distribution p comes from a family of probability distributions $p(\cdot; \theta)$ parameterized by $\theta \in \Theta$. When an efficient approximation algorithm exists, we describe a linear sketch from which a $(1 \pm \epsilon)$ -approximation $\hat{\ell}$ to $\ell(v)$ can be extracted with probability at least $2/3$. The form of the sketch is independent of the function g , so if $-\log p(x; \theta)$ satisfies our conditions for approximability for all $\theta \in \Theta$, we derive an approximation to each of the log-likelihoods $\ell(v; \theta)$ that are separately correct with probability $2/3$ each. If Θ is a discrete set then this allows us to find an approximate maximum likelihood estimator for θ with only an $O(\log |\Theta|)$ factor increase in storage. Recall, the maximum likelihood estimator is $\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \ell(\theta, v)$. Indeed, as usual we can drive the error probability down to any $\delta > 0$ with $O(\log 1/\delta)$ repetitions of the algorithm, so an additional factor of $O(\log |\Theta|)$ in the space complexity is enough. If $\hat{\ell}(\theta; v)$ denote our approximations, then $\hat{t} = \operatorname{argmin}_{\theta \in \Theta} \hat{\ell}(\theta; v)$ is the approximate maximum likelihood estimate, which has the guarantee that $\ell(\hat{t}; v) \leq (1 + \epsilon)\ell(\hat{\theta}; v)$. When Θ is not finite one may still be able to ap-

ply our results by discretizing Θ with $\operatorname{poly}(n)$ points and proceeding as above; whether this works depends on the behavior of $p(x; \theta)$ as a function of θ .

1.1.2 Utility Aggregates

Consider an online advertising service that charges its customers based on the number of clicks on the ad by web users. More clicks should result in a higher fee, but an exceptionally high number of clicks from one user make it likely that he is a bot or spammer. Customers may demand that the service discount the cost of these spam clicks, which means that the fee is a non-monotonic function of the number of clicks from each user.

The ads application is an example of a non-monotonic utility function. Many utility functions are naturally non-monotonic, and they can appear in streaming settings. For example, extremely low or high trading volume in a stock indicates abnormal market conditions. A low volume of network traffic might be a sign of damaged or malfunctioning equipment whereas a high volume of network traffic is one indicator of a denial-of-service attack.

1.1.3 Database Query Optimization

A prominent application for streaming algorithms, indeed one of the initial motivations [1], is to database query optimization. A query optimizer is a piece of software that uses heuristics and estimates query costs in order to intelligently organize the processing of a database query with the goal of reducing resource loads (time, space, network traffic, etc.). Scalable query optimizers are needed for scalable database processing [22, 16], and poor query efficiency has been cited as one of the main drawbacks of existing implementations of the MapReduce framework [19]. Streaming algorithms and sketches are a natural choice for such heuristics because of their low computational overhead and amenability to distributed processing. By expanding the character of statistics that can be computed efficiently over data streams, our results could allow database query optimizers more expressive power when creating heuristics and estimating query costs.

1.1.4 Encoding Higher Order Functions

An obvious generalization of the problem we consider is to replace the frequency vector f_i , $i \in [n]$, with a frequency matrix f_{ij} , $i \in [n]$ and $j \in [k]$, and ask whether there is a space efficient streaming algorithm that approximates $\sum_{i=1}^n g(f_{i,1}, \dots, f_{i,k})$. These could, for example, allow one to express more complicated database queries which first filter records based on one attribute and then sum up the function values applied to another attribute on the remaining records.

Let us demonstrate that when the coordinates of f are bounded and k is not too large, we can replace this sum with a function of a single variable. Suppose that $0 \leq f_{ij} \leq b-1 \in \mathbb{N}$, for all i, j . Upon receiving an update to coordinate (i, j) we replace it with b^j copies of i . The new frequencies are polynomially bounded if $b^k = \operatorname{poly}(n)$.

Call the new frequency vector $f' \in \mathbb{Z}^n$. The sequence of values f_{i1}, \dots, f_{ik} is immediately available as the base- b expansion of the number f'_i , so we are able to write $g'(f'_i) = g(f_{i1}, \dots, f_{ik})$, where g' first recovers the base- b expansion and then applies g . Our desire is to approximate $\sum_i g'(f'_i)$. Given even a well behaved function g , it is very unlikely that g' will be monotone, hence the need for an approximation

algorithm for non-monotonic functions if this approach is taken. It is also very likely that, because of the construction, g' has high local variability. This is a significant problem for algorithms that use only one pass over the stream (as we show by way of a lower bound), but we also present a two pass algorithm that is not sensitive to local variability.

1.2 Problem Definition

A stream of length m with domain $[n]$ is a list $D = ((i_1, \delta_1), (i_2, \delta_2), \dots, (i_m, \delta_m))$, where $i_j \in [n]$ and $\delta_j \in \mathbb{Z}$ for every $j \in [m]$. The *frequency vector* of a stream D is the vector $V(D) \in \mathbb{Z}^n$ with i^{th} coordinate $v_i = \sum_{j:i_j=i} \delta_j$. A processor can read the stream some number $p \geq 1$ times, in the order in which it is given, and is asked to compute a function of the stream. We allow the processor to use randomness and we only require that its output is correct with probability at least $2/3$.

Our algorithms are designed for the *turnstile streaming model*. In particular, there exists $M \in \mathbb{N}$ and it is promised that $V(D) \in \{-M, -M+1, \dots, M\}^n$ and the same holds for any prefix of the list D . Our lower bounds, on the other hand, hold in the more restrictive *insertion-only model* which has $\delta_j = 1$, for all j . We use $\mathcal{D}(n, m)$ to denote the set of all turnstile streams with domain $[n]$ and length at most m .

Given a function $g : \mathbb{R} \rightarrow \mathbb{R}$ and vector $V = (v_1, v_2, \dots, v_n)$, let

$$g(V) := \sum_{i=1}^n g(|v_i|).$$

DEFINITION 1. *Given $g : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$, $\epsilon > 0$, and $D \in \mathcal{D}(n, m)$ the problem of (g, ϵ) -SUM on stream D is to output an estimate \hat{G} of $g(V(D))$ such that*

$$P\left((1 - \epsilon)g(V(D)) \leq \hat{G} \leq (1 + \epsilon)g(V(D))\right) \geq 2/3.$$

We often omit ϵ and refer simply to g -SUM when the value of ϵ is clear from the context.

The choice of $2/3$ here is arbitrary, since given a g -SUM algorithm the success probability can be improved to $1 - 1/\text{poly}(n)$ by repeating it $O(\log n)$ times in parallel and taking the median outcome.

Our goal is to classify the space complexity of g -SUM. We ask: for which functions g can (g, ϵ) -SUM be solved in space that is sub-polynomial in n ? We call such functions *tractable*. An $\Omega(\epsilon^{-1})$ lower bound applies for nearly every function g , thus we only require that the algorithm solve (g, ϵ) -SUM whenever ϵ decreases sub-polynomially. Our main results separately answer this question in the case in which the processor is allowed one pass over the stream and in the case in which $O(1)$ passes are allowed. We are able to classify the space complexity of almost all functions, though a small set of functions remains poorly understood.

1.3 Organization

Section 2 gives a high-level description of our results and states our main theorems. We put off some of the more technical definitions until Section 3, which gives the definitions needed to precisely state our main results as well as background on the main techniques for our proofs. Next, Section 4 includes the intuition and proofs for our main theorems. It describes one and two-pass algorithms and lower

bounds. Section 5 discusses the functions that we are unable to characterize in more depth.

2. OUR RESULTS

To begin with, we separate the class of functions $g : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ into two complementary classes: “normal functions” and “nearly periodic functions”. In fact, we do this differently for 1-pass g -SUM than for multi-pass g -SUM. We define \mathbb{S} -normal functions, for studying one pass streaming space complexity, and \mathbb{P} -normal functions for multiple passes. The complements of these two sets are the \mathbb{S} -nearly periodic functions and \mathbb{P} -nearly periodic functions, respectively. The distinction between \mathbb{S} and \mathbb{P} is not important to understand our results at a high level, so we omit them for the rest of this section. See Section 5 for more about this class of functions. We postpone the technical definitions until Section 3, and give an informal overview here.

We are able to characterize the normal functions based on three properties which we call *slow-jumping*, *slow-dropping*, and *predictable*. Slow-jumping depends on the rate of increase. Roughly, a function is slow-jumping if it doesn’t grow much faster than $y = x^2$ at every scale. Slow-dropping governs the rate of decrease of a function. A function is slow dropping if it decreases no faster than sub-polynomially at every point. Finally, if a function is predictable then it satisfies a particular tradeoff between its growth rate and local variability. Interestingly, not just the properties themselves are important for our proofs, but so is the interplay between them. Precise definitions for these properties are given in Section 3.

Our proofs are by finding heavy hitters, for the upper bound, and by reductions from communication complexity for the lower bound. See Section 4.1 for more thorough discussion of how the three properties relate to heavy-hitters and the other techniques we use.

Zero-One Laws for Normal Functions

A nonnegative function f is called a sub-polynomial function if, for all $\alpha > 0$, $\lim_{x \rightarrow \infty} x^\alpha f(x) = \infty$ and $\lim_{x \rightarrow \infty} x^{-\alpha} f(x) = 0$. Formally, we study the class of functions g such that g -SUM can be solved with sub-polynomial accuracy $\epsilon = \epsilon(n)$ using only an amount of space which is sub-polynomial. We call such functions p -pass tractable if the algorithm uses p passes. In this paper we prove the following theorems.

THEOREM 2 (1-PASS ZERO-ONE LAW). *A function $g \in \mathcal{G}$ is 1-pass tractable and normal if and only if it is slow-jumping, slow-dropping, and predictable.*

THEOREM 3 (2-PASS ZERO-ONE LAW). *A function $g \in \mathcal{G}$ is 2-pass tractable and normal if and only if it is slow-dropping and slow-jumping.*

The main difference between the two theorems is that predictability is not needed in two passes. The message is that large local variability can rule out 1-pass approximation algorithms but not 2-pass approximation algorithms. Theorem 3 extends to any subpolynomial number of passes.

Most functions one encounters are normal, which include convex, concave, monotonic, polynomial, or trigonometric functions, as well as functions of regular variation and unbounded Lipschitz-continuous functions.

3. PRELIMINARIES

In describing the requirements and space efficiency of our algorithms for g -SUM, we use the set of sub-polynomial functions.

DEFINITION 4. A function $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is sub-polynomial if for any $\alpha > 0$, the following two conditions are satisfied: a) $\lim_{x \rightarrow \infty} x^\alpha f(x) = \infty$ and b) $\lim_{x \rightarrow \infty} x^{-\alpha} f(x) = 0$.

We use $\text{subpoly}(x)$ to represent the set of sub-polynomial functions in the variable x . Examples of sub-polynomial functions include polylogarithmic functions and some with faster growth, like $2^{\sqrt{\log n}}$. The set of poly-logarithmic functions is a subset of the sub-polynomial functions. Formally, we study the class of functions g such that (g, ϵ) -SUM can be solved with any accuracy $\epsilon \in \text{subpoly}(n)$, using only sub-polynomial space. Our algorithms assume an oracle for computing g and that the storage required for the value $g(x)$ is sub-polynomial in x . Regardless, our sketches are sub-polynomial in size, but without these assumptions it could take more than sub-polynomial space to compute the approximation from the sketch. We restrict ourselves to the case $M \in \text{poly}(n)$. This restriction is reasonable because if it grows, say, exponentially in n , then an algorithm can store the entire frequency vector and compute $g(V(D))$ exactly in $\text{polylog}(nM)$ space.

DEFINITION 5. A function $g : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ is p -pass tractable if given any sub-polynomial function $h(x)$ and any $\epsilon \geq 1/h(nM)$ there exists a sub-polynomial function $h^*(x)$ and a p -pass algorithm \mathcal{A} that solves (g, ϵ) -SUM for all streams in $\mathcal{D}(n, m)$ and $n, M \geq 1$ using no more than $h^*(nM)$ space in the worst case.

Replacing the storage requirement, which is $h^*(nM)$, with $h^*(n) \log M$ would also be natural, but since we consider $M \in \text{poly}(n)$ the two definitions are equivalent. It simplifies our notation a little bit to write $h^*(nM)$. Unless otherwise noted, for the rest of this paper we require that $g(0) = 0$ and $g(x) > 0$, for all $x > 0$. The choice $g(0) = 0$ is equivalent to requiring that $g(V(D))$ does not depend on the particular choice of n in the model; specifically it avoids the following behavior: if $g(0) \neq 0$ then given a single stream D the value of $g(V(D))$ differs depending on the choice of the dimension, even though the stream remains the same. Functions with $g(0) \neq 0$ may be of interest for some applications. The laws for these functions are very similar to the case when $g(0) = 0$ and we provide them in Appendix A.

As shown in [10], functions with $g(x) = 0$, for some integer $x > 0$, are not n^α -pass tractable for any $\alpha < 1$, unless g is nonnegative and periodic with period $\min\{x > 0 \mid g(x) = 0\}$, which must exist if g is periodic and $g(0) = 0$. It is also shown in [10] that if a non-linear function g takes both positive and negative values, then g is not n^α -pass tractable, for any $\alpha < 1$. Without loss of generality, we can assume $g(1) = 1$ because a multiplicative approximation algorithm for the function $g'(x) := g(x)/g(1)$ is also a multiplicative approximation algorithm for g . Finally, for simplicity of notation we will often extend the domain of g symmetrically to \mathbb{Z} , i.e., setting $g(x) = g(-x)$, for all $x \in \mathbb{Z}_{\geq 0}$, which allows us the simpler notation $g(|v_i|) = g(v_i)$. In summary, we study the functions in the class

$$\mathcal{G} \equiv \{g : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}, g(0) = 0, g(1) = 1, \forall x > 0, g(x) > 0\}$$

Our results are based on three characterizations of the variation of the functions. We now state the technical definitions mentioned in Section 2.

DEFINITION 6. A function $g \in \mathcal{G}$ is slow-jumping if for any $\alpha > 0$, there exists $N > 0$, such that for any $x < y$ if $y \geq N$, then $g(y) \leq \lfloor \frac{y}{x} \rfloor^{2+\alpha} x^\alpha g(x)$.

Slow-jumping is a characterization of the rate of growth of a function. Examples of slow-jumping functions are x^p , for $p \leq 2$, $x^2 2^{\sqrt{\log x}}$, and $(2 + \sin x)x^2$, while functions like 2^x and x^p , for any $p > 2$, are not slow jumping because they grow too quickly.

DEFINITION 7. A function $g \in \mathcal{G}$ is slow-dropping if for any $\alpha > 0$ there exists $N > 0$, such that for any $x < y$, if $y \geq N$, then $g(y) \geq g(x)/y^\alpha$.

Whether a function is slow-dropping is determined by its rate of decrease. The functions $(\log_2 1 + x)^{-1} \mathbf{1}(x > 0)$ and $(2 + \sin x)x^2$ are slow-dropping, but any function with polynomial decay, i.e. x^{-p} , for $p > 0$, is not slow-dropping.

Given a function $g, x \in \mathbb{N}$, and $\epsilon > 0$ define the set: $\delta_\epsilon(g, x) := \{y \in \mathbb{N} \mid |g(y) - g(x)| \leq \epsilon g(x)\}$.

DEFINITION 8. A function g is predictable if for every $0 < \gamma < 1$ and sub-polynomial ϵ , there exists N such that for all $x \geq N$ and $y \in [1, x^{1-\gamma})$ with $x + y \notin \delta_{\epsilon(x)}(g, x)$ we have $g(y) \geq x^{-\gamma} g(x)$.

Predictability is governed by the local variability of a function and its rate of growth. For example, the function $g(x) = x^2$ is predictable because $g(x+y)/g(x) = (1 + \frac{y}{x})^2 \approx 1$ when $y \ll x$, or more precisely $\pi_\epsilon(x) \supseteq \{y \mid |x - y| \leq \epsilon x/3\}$. That function has low local variability. On the other hand, the function $(2 + \sin x)\mathbf{1}(x > 0)$ is locally highly variable but still predictable. Notice that even $x + 1 \notin \pi_\epsilon(x)$, for say $\epsilon = 0.01$, but $g(1)/g(x) \geq 1/3$ so the inequality in the definition of predictability is satisfied. A negative example is the function $(2 + \sin x)x^2$, which is not predictable because it varies quickly (by a multiplicative factor of 3) and grows with x .

As we will show, the three conditions above can be used to characterize nearly every function in \mathcal{G} in both the single-pass and $O(1)$ -pass settings. We remark here that the characterization has already been completed for monotonic functions. The tractability condition for nondecreasing g proved by [6] is equivalent to g being slow-jumping and predictable. For nonincreasing functions, it is a consequence of [5] that polynomially decreasing functions are not tractable while sub-polynomially decreasing functions are tractable. The “nearly periodic” functions are formally defined next.

Let \mathbb{S} be the set of non-increasing sub-polynomial functions on domain $\mathbb{Z}_{\geq 0}$ and \mathbb{P} be the set of strictly increasing polynomial functions on $\mathbb{Z}_{\geq 0}$. We now define the sets of \mathbb{S} -nearly periodic functions and \mathbb{P} -nearly periodic functions. The motivation for these definitions jumps ahead to our lower bounds, but we will explain it here. The reductions used in some of our lower bounds boil down to finding three integers $x < y$ and $x + y$ such that $g(x) \gg g(y)$ and $g(x) \not\approx g(x + y)$. The reductions fail when $g(x) \approx g(x + y)$, where the meaning of “ \approx ” depends on whether we are bounding 1-pass or multi-pass algorithms. For the 1-pass reduction

¹ $\mathbf{1}(\cdot)$ denotes the indicator function.

to fail it is necessary that $\frac{1}{g(x)}|g(x) - g(x+y)|$ decreases as x increases, hence the \mathbb{S} -nearly periodic functions. The 2-pass reduction fails as long as $\max(\frac{g(x)}{g(x+y)}, \frac{g(x+y)}{g(x)})$ is not polynomially large in y , hence the \mathbb{P} -nearly periodic functions.

DEFINITION 9. *Given a set of functions \mathcal{S} , call $g(x)$ \mathcal{S} -nearly periodic, if the following two conditions are satisfied.*

1. *There exists $\alpha > 0$ such that for any constant $N > 0$ there exists $x, y \in \mathbb{N}$, $x < y$ and $y \geq N$ such that $g(y) \leq g(x)/y^\alpha$. Call such a y an α -period of g ;*
2. *For any $\alpha > 0$ and any error function $h \in \mathcal{S}$ there exists $N_1 > 0$ such that for all α -periods $y \geq N_1$ and all $x < y$ such that $g(y)y^\alpha \leq g(x)$, we have $|g(x+y) - g(x)| \leq \min\{g(x), g(x+y)\}h(y)$.*

A function g is \mathcal{S} -normal if it is not \mathcal{S} -nearly periodic.

The first condition states that the function is not slow-dropping. For example, if g is bounded then there is an increasing subsequence along which g converges to 0 polynomially fast. To understand the second condition, take h to be a decreasing function or a small constant. In loose terms, it states that if $x < y$ and $g(x) \gg g(y)$, then $g(x) \approx g(x+y)$. The choice of the set of functions \mathcal{S} determines the relative accuracy implied by “ \approx ”. We will apply the definition with \mathcal{S} set to either \mathbb{S} or \mathbb{P} , that is, with either subpolynomially or polynomially small relative accuracy.

An \mathbb{S} -nearly periodic function (it is also \mathbb{P} -nearly periodic) can be constructed as follows. For each $x \in \mathbb{N}$, let $i_x = \max\{j : 2^j | x\}$ and let $g_{np}(0) = 0$ and $g_{np}(x) = 2^{-i_x}$. For example, $g_{np}(1) = 1$, $g_{np}(2) = 1/2$, $g_{np}(3) = 1$, $g_{np}(4) = 1/4$, etc. The proof that g is nearly periodic appears in Section 5, but for now notice that it satisfies the first part of the definition because $g_{np}(2^k) = 2^{-k}$ and, as an example of the second part, we have $g_{np}(2^k + 1) = g_{np}(1) = 1$. One may guess that a streaming algorithm for such an erratic function requires a lot of storage. In fact, g_{np} can be approximated in polylogarithmic space! Section 5 has the proof.

The following containment follows directly from the definition of the nearly periodic functions.

PROPOSITION 10. *Every \mathbb{S} -nearly periodic function is also \mathbb{P} -nearly periodic. Every \mathbb{P} -normal function is \mathbb{S} -normal.*

3.1 Heavy Hitters, Communication Complexity, and CountSketch

Our sub-polynomial space algorithm is based on the Recursive Sketch of Braverman and Ostrovsky [7], which reduces the problem to that of finding heavy hitters.

DEFINITION 11. *Given a stream $D \in \mathcal{D}(n, m)$ and a function g , call $j \in [n]$ a (g, λ) -heavy hitter of $V(D)$ if $g(|v_j|) \geq \lambda \sum_{i \neq j} g(|v_i|)$. We will use the terminology λ -heavy hitter when g is clear from the context.*

Given a heaviness parameter λ and approximation accuracy ϵ , a heavy hitters algorithm returns a set $\{i_1, i_2, \dots\}$ containing every (g, λ) -heavy hitter and also returns a $(1 \pm \epsilon)$ -approximations to $g(v_{i_j})$ for each i_j in the set. Such a collection of pairs is called a (g, λ, ϵ) -cover.

DEFINITION 12. *Given a stream $D \in \mathcal{D}(n, m)$, a (g, λ, ϵ) -cover of $V(D)$ is a set of pairs $\{(i_1, w_1), (i_2, w_2), \dots, (i_t, w_t)\}$, where $t \leq n$ and the following are true:*

1. $\forall k \in [t]$, $|w_k - g(|v_{i_k}|)| \leq \epsilon g(|v_{i_k}|)$, and
2. *if j is a (g, λ) -heavy hitter of $V(D)$, then there exists $k \in [t]$ such that $i_k = j$.*

Formally, a p -pass $(g, \lambda, \epsilon, \delta)$ -heavy hitter algorithm is a p -pass streaming algorithm that outputs a (g, λ, ϵ) -cover of $V(D)$ with probability at least $(1 - \delta)$.

THEOREM 13 ([7]). *Let $\lambda = \frac{\epsilon^2}{\log^3 n}$ and $\delta = \frac{1}{\log n}$. If there exists a $(g, \lambda, \epsilon, \delta)$ -heavy hitter algorithm using s bits of space, then there exists a (g, ϵ) -SUM algorithm using $O(s \log n)$ bits of space.*

A similar approach has been used in other streaming algorithms; it was pioneered by Indyk and Woodruff [15] for the problem of approximating the frequency moments. Their algorithm uses the streaming algorithm CountSketch of Charikar, Chen, and Farach-Colton [9], and ours will as well. Given λ , ϵ , and δ , a CountSketch using space $O(\frac{1}{\lambda \epsilon^2} \log \frac{n}{\delta} \log M)$ gives an approximation \hat{v}_i to the frequency of every item $i \in [n]$ with the following guarantee: with probability at least $(1 - \delta)$, $|v_i - \hat{v}_i| \leq \epsilon \sqrt{\lambda F_2}$ for all $i \in [n]$. By adjusting the parameters slightly we can treat the output of CountSketch($\lambda, \epsilon, \delta$) as a set of $k = O(1/\lambda)$ pairs (i_j, \hat{v}_{i_j}) such that $\{i_j\}_{j=1}^k$ contains all of the λ -heavy hitters for F_2 and $|v_{i_j} - \hat{v}_{i_j}| \leq \epsilon \left(\sum_{j > k} \bar{v}_j^2\right)^{1/2} \leq \epsilon \sqrt{\lambda F_2}$, for all $j = 1, 2, \dots, k$, where $(\bar{v}_1, \dots, \bar{v}_n)$ is a reordering of the frequency vector from largest to smallest absolute magnitude.

We establish our lower bounds using communication and information complexity. Several of our lower bounds can be shown by reduction from standard INDEX and DISJ problems, as well as a combination of them (related combined problems were used in [10, 21]). In the INDEX(n) problem, there are two players Alice, given a set $A \subseteq [n]$, and Bob, given $b \in [n]$. Alice sends a message to Bob, and with probability at least $2/3$, Bob must determine whether $b \in A$. Any randomized one-way protocol for INDEX(n) requires Alice to send $\Omega(n)$ bits [18]. In an instance of DISJ(n, t) there are t players each receiving a subset of $[n]$ with the promise that either the sets are pairwise disjoint, or there is a single item that every set contains and other than this single item, the sets are pairwise disjoint. The players must determine which case they are in. The randomized, unrestricted communication complexity of DISJ(n, t) is $\Omega(n/t)$ [4, 8, 11, 17]. In the DISJ+IND(n, t) problem, $t + 1$ players are given sets $A_1, A_2, \dots, A_{t+1} \subseteq [n]$, such that $|A_{t+1}| = 1$, with the promise that either the sets are disjoint or there is a single element contained in every set and they are otherwise disjoint. A standard argument shows the one-way communication complexity of DISJ+IND(n, t) is $\Omega(n/t(\log n))$ (see Theorem 47).

To obtain stronger lower bounds for the nearly periodic function sums, we use the information complexity framework.

DEFINITION 14 (ShortLinearCombination PROBLEM). *Let $u = (u_1, u_2, \dots, u_r)$ be a vector in \mathbb{Z}^r for some integer r . Let $d > 0$ be an integer not in $\{|u_1|, |u_2|, \dots, |u_r|\}$. A stream S with frequency vector v is given to the player, v is promised to be from $V_0 = \{u_1, u_2, \dots, u_r, 0\}^n$ or $V_1 = \{EMB(v, i, e) \mid v \in V_0, i \in [n], e \in \{-d, d\}\}$, where $EMB(v, i, e)$ is to embed e onto the i -th coordinate. The (u, d) -DIST problem is to distinguish whether $v \in V_0$ or $v \in V_1$.*

This problem is a generalization of the DISJ(n, t) problem. We show a $\Omega(n/q^2)$ bits lower bound for it where $q = \sum_{i=1}^r |q_i|$ and q_i are the integers such that $\min_q \{q_1, q_2, \dots, q_r \mid \sum_{i=1}^r q_i u_i = d\}$. For more details, we refer the reader to Appendix C.

4. ZERO-ONE LAWS FOR NORMAL FUNCTIONS

First, we develop some intuition about why the three conditions, slow-dropping, slow-jumping, and predictable, characterize tractable normal functions. Next comes the two pass and one pass algorithms followed by the lower bounds. Finally, we prove Theorem 2 and 3 in Section 4.6.

4.1 Our Techniques

The upper and lower bounds are both established using heavy hitters. Using a $\frac{\epsilon^2}{\log^3 n}$ -heavy hitters subroutine that identifies each heavy hitter H and also approximates $g(f_H)$, the algorithm of Braverman and Ostrovsky [7] (Theorem 13) solves g -SUM with $O(\log n)$ storage overhead. We will show that if g is tractable then a subpolynomially sized CountSketch suffices to find heavy hitters for g . On the lower bounds side, typical streaming lower bounds, going back to the work of Alon, Matias, and Szegedy, are derived from reductions to communication complexity that only require the streaming algorithm to detect the presence of a heavy hitter. Thus the problem of characterizing tractable functions can be reduced to characterizing the set of functions admitting subpolynomial heavy hitters algorithms.

Now let us explain how each of the three properties relates to heavy hitters. For this discussion let $\alpha = \epsilon^2 / \log^3 n$ be the heaviness parameter needed by Braverman's and Ostrovsky's algorithm and let $H \in [n]$ be the identity of some α -heavy hitter.

Slow-jumping and Slow-dropping.

A function that satisfies these two conditions cannot increase much faster than a quadratic function nor can it decrease rapidly. This means that if H is a heavy hitter for g then it is also a heavy hitter for F_2 . Concretely, suppose that there exists an increasing subpolynomial function h such that $g(y)/g(x) \leq (y/x)^2 h(y)$ and $g(x) \geq g(y)/h(y)$, for all $x \leq y$, so that g is slow-jumping and slow-dropping (Propositions 15 and 16 show that such an h always exists if g is slow-jumping and slow-dropping). A bit of algebra shows that $g(f_H) \geq \alpha \sum_i g(f_i)$ implies $f_H^2 \geq \frac{\alpha}{h(M)^2} \sum_j f_j^2$, so H is a $\alpha' := \frac{\alpha}{h(M)^2}$ -heavy hitter for F_2 . CountSketch suffices to identify all α' -heavy hitters for F_2 with $O(\alpha' \log^2 n)$ bits, so this gives an α -heavy hitters algorithm for g that still uses subpolynomial space.

If a function is not slow-jumping then a heavy hitter may be too small to detect and we prove a lower bound by reduction from DISJ+IND in a very similar fashion to the original lower bound for approximating the frequency moments by Alon, Matias, and Szegedy.

If a function is not slow-dropping then f_H may be hidden below many small frequencies. In this case, a reduction from the communication complexity of INDEX usually works to establish a one-pass lower bound (two player DISJ can be used for a multipass lower bound). Here is a preview of the reduction that will also explain the genesis of

the nearly periodic functions. Suppose $g(1) = 1$ and there is an increasing sequence $n_k \in \mathbb{N}$ such that $g(n_k) \leq 1/n_k$. Alice and Bob can solve an instance ($A \subseteq [n_k], b \in [n_k]$) of INDEX as follows. Alice creates a stream where every item in A has frequency n_k and the items in $[n_k] \setminus A$ do not appear. To Alice's stream Bob adds one copy of his item and they jointly run a streaming algorithm to approximate $g(f)$. The result is an approximation to either $|A|g(n_k) + 1$ or $(|A| - 1)g(n_k) + g(n_k + 1)$. If $g(n_k + 1)$ differs significantly from $g(n_k) + 1 \approx 1$ then Bob can determine whether $b \in A$ by checking the approximation to $g(f)$. The reduction fails if $g(n_k + 1) \approx 1$ and such a function may be tractable! A prime example is g_{np} from Section 3, which we demonstrate in Section 5.

At this point, one might have in mind our two pass heavy hitters algorithm. The algorithm is to use a CountSketch on the first pass to identify a subpolynomial size set of items containing all α -heavy hitters and then use the second pass just to tabulate the frequency of each of those items exactly. The details and proof of correctness for that algorithm are in Section 4.2. Local variability of the function is irrelevant for two passes because we can compute the frequencies exactly during the second pass, and that is why a function need not satisfy predictability to be two pass tractable. A one pass algorithm will have to get by with approximate frequencies.

Predictable.

Predictability is a condition on the local variability of the function. It roughly says that if $y \ll x$ then either $g(x + y) \approx g(x)$ or $g(y)$ is on the same scale² as $g(x)$. The first conclusion has an obvious interpretation; it means that substituting an approximation to x yields an approximation to $g(x)$. CountSketch returns an approximation to the frequency of every heavy hitter, so, in the first case, we can just as well use the approximation and there is no need to determine the frequency exactly.

The second conclusion allows local variability, but it must be accompanied by some "global" property of g , namely, g does not change a whole lot over the interval $[y, x]$. For an illustration of how predictability works, suppose that H is a heavy hitter and there is some $y \ll f_H$ such that $g(y) \geq 2g(f_H)$. Our algorithm cannot substitute $f_H + y$ for f_H because this could lead to too much error. But, approximating f_H to better accuracy than $\pm y$ using a CountSketch, one would generally need more than $(f_H/y)^2$ counters, which could be very large. Notice that y , if it occurs as a frequency in the stream and suppose it does, is a heavy hitter, and if g is slow-jumping and slow-dropping then y is also a heavy hitter for F_2 by our previous argument. It follows that the error in the frequency estimate derived from the CountSketch is less than y , in particular an estimate of f_H derived from the CountSketch will be more accurate than $\pm y$. The result is that we get an error estimate for the CountSketch that is improved over a naive analysis. The approximate frequency \hat{f}_H satisfies $g(\hat{f}_H) \approx g(f_H)$ and is accurate enough for Braverman and Ostrovsky's algorithm. When g is locally variable but not predictable, we get a one pass lower bound by a reduction from INDEX.

²More precisely, $g(y)$ is not much smaller than $g(x)$, if g is slow-dropping then $g(y)$ cannot be much larger than $g(x)$, either.

4.2 Two Pass Algorithm

We will now show that any λ -heavy hitter for a slow-dropping, slow-jumping, \mathbb{S} -normal function is also F_2 λ' -heavy for λ' modestly smaller than λ . This means that we can use the CountSketch to identify heavy hitters in one pass for the Recursive Sketch. For a heavy hitters algorithm we also need a $(1 \pm \epsilon)$ -approximation to each item's contribution to g -SUM, we can easily accomplish this in the second pass by exactly tabulating the frequency of each heavy hitter identified in the first pass. This algorithm works as long as we identify every heavy hitter in the first pass, which CountSketch guarantees, and the space required remains small provided that we do not misclassify too many non-heavy hitters (so that we do not tabulate too many values in the second pass). The procedure is formally defined in Algorithm 1. Note that by Proposition 10 \mathbb{P} -normal functions are also \mathbb{S} -normal, hence, the algorithm works for slow-dropping and slow-jumping \mathbb{P} -normal functions.

The next two propositions describe equivalent definitions of slow-dropping and slow-jumping that are useful in describing the algorithm.

PROPOSITION 15. *$g \in \mathcal{G}$ is slow-dropping if and only if there exists a sub-polynomial function h such that for all $y \in \mathbb{N}$ and $x < y$ we have $g(x) \leq g(y)h(y)$.*

PROOF. The “if” direction follows immediately from the definition of the class of sub-polynomial functions.

For the “only if”, suppose that $g(x)$ is slow-dropping. Specifically, for any $\alpha > 0$ there exists $N > 0$ such that for all $y \geq N$ and $x < y$ we have $y^\alpha g(y) \geq g(x)$. Let N_i be the least integer such that $y^{1/i} g(y) \geq g(x)$, for all $y > N_i$ and $x < y$. The sequence N_i is nondecreasing. Consider the increasing function $i(y) = \max\{i : N_i \leq y\}$ and set $h(y) = y^{1/i(y)}$. For all $x < y$ we have $h(y)g(y) \geq g(x)$ by construction.

To see that h is sub-polynomial, let $\beta > 0$. First, $h \geq 1$ hence $h(y)y^\beta \rightarrow \infty$. Second, let $i^* \geq 2/\beta$ then for all $y \geq N_{i^*}$ we have

$$h(y) \leq y^{1/i^*} \leq y^{\beta/2}.$$

Thus $h(y)y^{-\beta} \rightarrow 0$, which completes the proof. \square

PROPOSITION 16. *$g \in \mathcal{G}$ is slow-jumping if and only if there exists a sub-polynomial function $h(x)$ such that for any $x < y$ we have $g(y) \leq \lfloor y/x \rfloor^2 h(\lfloor y/x \rfloor x)g(x)$.*

PROOF. Again, the “if” direction follows immediately from the definition of the class of sub-polynomial function. The reverse direction follows from the same argument as Proposition 15. \square

Given g and a sub-polynomial accuracy ϵ , Propositions 15 and 16 each imply the existence of a non-decreasing sub-polynomial function. By taking the point-wise maximum, we can assume that these are the same sub-polynomial function $H : \mathbb{N} \rightarrow \mathbb{R}$. In particular H satisfies:

- $g(y) \geq g(x)/H(y)$, for all $x < y$, and
- $g(y) \leq (y/x)^2 y^\alpha H(y)g(x)$, for all $x < y$.

The space used by our algorithm depends on the sub-polynomial functions governing the slow-dropping and slow-jumping of g . If these are polylogarithmic then the algorithm is limited to polylogarithmic space.

LEMMA 17. *Let g be a function that is slow-jumping and slow-dropping. There exists a sub-polynomial function h such that for any $D \in \mathcal{D}(n, m)$ with frequencies v_1, v_2, \dots, v_n , if $g(v_i) \geq \lambda \sum_j g(v_j)$ then*

$$v_i^2 \geq \frac{\lambda}{h(\lfloor v_i \rfloor)} \sum_{|v_j| < |v_i|} v_j^2.$$

PROOF. By Proposition 16, there exists a nondecreasing sub-polynomial function h such that for all $y \in \mathbb{N}$ and $x < y$ we have $g(y) \leq (y/x)^2 h(y)g(x)$.

For any j that satisfies $|v_j| < |v_i|$, we have $g(v_i) \leq g(v_j)(\frac{v_i}{v_j})^2 h(\lfloor v_i \rfloor)$. Therefore,

$$g(v_i) \geq \lambda \sum_j g(v_j) \geq \lambda \sum_{|v_j| < |v_i|} \frac{g(v_i)}{h(\lfloor v_i \rfloor)} \left(\frac{v_j}{v_i}\right)^2.$$

By rearranging, we find

$$v_i^2 \geq \frac{\lambda}{h(\lfloor v_i \rfloor)} \sum_{|v_j| < |v_i|} v_j^2 \geq \frac{\lambda}{h(\lfloor v_i \rfloor)} \sum_{|v_j| < |v_i|} v_j^2.$$

\square

Algorithm 1 A 2-pass $(g, \lambda, 0, \delta)$ -heavy hitters algorithm.

- 1: **procedure** 2-PASS HEAVY HITTERS($g, \lambda, \epsilon, \delta$)
 - 2: **First Pass:**
 - 3: $S \leftarrow \text{CountSketch}(\frac{\lambda}{2H(M)}, \frac{1}{3}, \delta)$ discarding the frequency estimates
 - 4: **Second Pass:**
 - 5: Compute v_j for all $i \in S$
 - 6: **return** (j, v_j) , for all $j \in S$
-

Algorithm 1 computes CountSketch³ with $r = O(\log n)$ and $b = O(\frac{H(M)}{\lambda \epsilon^2})$ over the stream and extracts the $2H(M)/\lambda$ items with the highest estimated frequencies. Denote this list as $(i_j, \hat{v}_{i_j})_{j=1}^{2H(M)/\lambda}$. The 2-pass algorithm then discards the estimated frequencies \hat{v}_{i_j} and tabulates the true frequency of each item i_j on the second pass.

LEMMA 18. *If g is a function that is slow-dropping and slow-jumping then the CountSketch used by Algorithm 1 finds all of the (g, λ) -heavy hitters.*

PROOF. From the slow-dropping condition and Lemma 17, for every $x < y \in \mathbb{N}$ we have

1. $g(y)H(y) \geq g(x)$ and
2. $g(v_i) \geq \lambda \sum_j g(v_j)$ implies

$$v_i^2 \geq \frac{\lambda}{H(M)} \sum_{|v_j| < |v_i|} v_j^2.$$

Let $D \in \mathcal{D}(n, m)$, suppose i satisfies $g(v_i) \geq \lambda \sum_j g(v_j)$. Since

$$g(v_i) \geq \lambda \sum_{|v_j| \geq |v_i|} g(v_j) \geq \frac{\lambda}{H(M)} \sum_{|v_j| \geq |v_i|} g(v_i),$$

there are at most $H(M)\lambda^{-1}$ items with frequencies as large or larger than v_i in magnitude.

³A CountSketch is a matrix with r and b rows. See [9] for details about these parameters.

Thus, a CountSketch for $\lambda/2H(M)$ -heavy hitters, as is used by Algorithm 1, serves to identify the λ -heavy heavy hitters for g in one pass. \square

THEOREM 19. *A function $g \in \mathcal{G}$ is 2-pass tractable and \mathbb{S} -normal if it is slow-dropping and slow-jumping.*

PROOF. If a function is slow-dropping then it is \mathbb{S} -normal. By Theorem 13, it is sufficient to show that Algorithm 1 is a sub-polynomial memory $(g, \lambda, \epsilon, \delta)$ -heavy hitters algorithm for $\lambda = \frac{\epsilon^2}{\log^3 n}$, $\epsilon = \frac{1}{H(Mn)}$, and $\delta \leq \frac{1}{\log n}$.

This is a 2-pass $(g, \lambda, 0, \delta)$ -heavy hitters algorithm, and it requires $O(\frac{h(M)}{\epsilon^2} \log^4 n \log n M \log \log n)$ space. The algorithm is correct with probability at least $(1 - \delta)$ because this is true for the CountSketch. Therefore, g is a 2-pass tractable function. \square

4.3 One Pass Algorithm

The only impediment to reducing the two pass algorithm to a single pass is local variability of the function. The algorithm must simultaneously identify heavy hitters and estimate their frequencies, but local variability could require tighter estimates than seem to be available from a sub-polynomial space CountSketch data structure. If an \mathbb{S} -normal function is predictable then we almost immediately overcome this barrier. In particular, the definition means that given any point x and a small error $y = o(x)$ either $|g(x+y) - g(x)| \leq \epsilon g(x)$ or $g(y)$ is reasonably large, and while it is clear how the first case helps us with the approximation algorithm, the second is less clear. The slow-dropping and slow-jumping conditions play a role, as we now explain.

PROPOSITION 20. *A function g is predictable if and only if for every sub-polynomial $\epsilon > 0$ there exists a sub-polynomial function h such that for all $x \in \mathbb{N}$ and $y \in [1, x/h(x))$ satisfying $x + y \notin \delta_{\epsilon(x)}(g, x)$ it holds that $g(y) \geq g(x)/h(x)$.*

PROOF. The “if” direction follows directly from the definition of sub-polynomial functions.

To prove the “only if” direction, let $\gamma_i = 1/i$. Since g is predictable, there exists N_i such that either $x + y \in \delta_{\epsilon(x)}(g, x)$ or $g(y) \geq x^{-\gamma} g(x)$ for $x \geq N_i$ and $y \in [1, x^{1-\gamma})$. Let $i(z) = \max\{i \in \mathbb{N} \mid N_i \leq z\}$, for $z \geq N_1$, and define $h(z) = z^{1/i(z)}$, for $z \geq N_1$, and $h(z) = \max_{1 \leq x', y' < N_1} g(x')/g(y')$, for $z < N_1$. Then h is a sub-polynomial function. The claim is that h satisfies the conclusion of the lemma. Let $x \in \mathbb{N}$ and $y \in [1, x/h(x))$. If $x < N_1$ then $y \leq x < N_1$ because $h \geq 1$. Furthermore,

$$g(y)h(x) = g(y) \max_{1 \leq x', y' < N_1} \frac{g(x')}{g(y')} \geq g(y) \frac{g(x)}{g(y)} = g(x).$$

If $x \geq N_1$, let $i = i(x) \geq 1$, so that $x \geq N_i$ and $h(x) = x^{1/i}$. Since g is predictable, if $y \in [1, x/h(x))$ and $y + x \notin \delta_{\epsilon(x)}(g, x)$ then $g(y) \geq x^{-1/i} g(x) = g(x)/h(x)$. \square

Let $r_\epsilon(x) := \max\{y \mid x + y' \in \delta_\epsilon(x), \text{ for all } |y'| \leq y\}$.

LEMMA 21 (PREDICTABILITY BOOSTER). *If g is slow-dropping, slow-jumping, and predictable, then there is a sub-polynomial function h such that for all $y \in [r_\epsilon(x)+1, x/h(x))$ we have $g(y) \geq g(x)/h(x)$.*

PROOF. According to Proposition 20, there exists a sub-polynomial function h' such that for all $y \in [1, x/h'(x))$ if $x + y \notin \delta_{\epsilon/2}(x)$ then $g(y) \geq g(x)/h'(x)$. Without loss of

generality h' governs the slow-jumping of g as well, hence $g(x) \leq (x/x')^2 h'(x) g(x')$ for all $x' \leq x$.

We consider two cases:

1. $|g(x) - g(x + r_\epsilon(x) + 1)| > \epsilon g(x)$ and
2. $|g(x) - g(x - r_\epsilon(x) - 1)| > \epsilon g(x)$.

In the first case, by the definition of r_ϵ and predictability we have $g(r_\epsilon(x) + 1) \geq g(x)/h'(x)$. In the second case, let $x' = x - r_\epsilon - 1$. If $g(x') \geq 2g(x)$ then $|g(x') - g(x)| > \epsilon g(x')$, as long as $\epsilon < 1/2$. Otherwise, $|g(x') - g(x)| \geq \epsilon g(x) \geq \frac{\epsilon}{2} g(x')$. Now, from predictability at x' and the definition of h' , we find that

$$g(r_\epsilon(x) + 1) = g(x - x') \geq \frac{g(x')}{h'(x')} \geq \frac{g(x)}{4h'(x)^2},$$

where the last inequality follows because g is slow-jumping and $x' \geq x/2$ (w.l.o.g., choose $h'(x) > 2$).

Now, since g satisfies the slow dropping condition there exists a nondecreasing sub-polynomial function h'' such that $g(y') \geq g(r')/h''(y')$ for all $y' \in \mathbb{N}$ and $r' \leq y'$. Thus, if $y \in [r_\epsilon(x) + 1, x/h(x))$

$$g(y) \geq \frac{g(r_\epsilon(x) + 1)}{h''(y)} \geq \frac{g(x)}{4h'(x)^2 h''(x)},$$

so we take h to be the product of $4(h')^2$ and h'' . \square

Recall the sub-polynomial function H from the last section and let it also satisfy Lemma 21 with ϵ replaced by $\epsilon/2$. In particular H now satisfies three conditions:

- $g(y) \geq g(x)/H(y)$, for all $x < y$, and
- $g(y) \leq (y/x)^2 H(y)g(x)$, for all $x < y$.
- for all $y \in [r_{\epsilon/2}(x)+1, x/h(x))$ we have $g(y) \geq g(x)/H(x)$.

Let y be the (additive) error in our estimate of frequency x of a (g, λ) -heavy hitter. If $r_\epsilon(x) > |y|$ then we are fine. If not, consider the point $x + r_\epsilon(x) + 1$ (or $x - r_\epsilon(x) - 1$), which has $|g(x + r_\epsilon(x) + 1) - g(x)| > \epsilon g(x)$ and therefore $g(r_\epsilon(x) + 1) > g(x)/H(M)$. Now, since x is the frequency of a (g, λ) -heavy hitter it happens that $r_\epsilon(x) + 1$, were it a frequency in the stream, would be $(g, \lambda/H(M))$ -heavy and thus $\lambda/H(M)^2$ -heavy for F_2 . Presuming that there are not too many items in the stream with frequency larger than $r_\epsilon(x)$, this implies that CountSketch produces an estimate for x with error smaller than $r_\epsilon(x)$. Now, since g satisfies the slow dropping condition there cannot be too many frequencies larger than $r_\epsilon(x)$ in the stream because

$$g(z) \geq \frac{g(r_\epsilon(x) + 1)}{H(M)} \geq \frac{g(x)}{H(M)^2}, \text{ for all } z > r_\epsilon,$$

which implies that there are at most $\frac{1}{\lambda} H(M)^2$ frequencies in this range.

4.4 One Pass Lower Bounds

The proof of the following theorem is broken up into Lemmas 23, 24 and 25.

THEOREM 22. *If $g \in \mathcal{G}$ is a 1-pass tractable \mathbb{S} -normal function, then g is slow-jumping, slow-dropping, and predictable.*

PROOF. Immediate from Lemmas 23, 24, and 25. \square

Algorithm 2 A 1-pass $(g, \lambda, \epsilon, \delta)$ -heavy hitters algorithm.

1: **procedure** 1-PASS HEAVY HITTERS($g, \lambda, \epsilon, \delta$)
2: $\hat{S}, \hat{V} \leftarrow \text{CountSketch}(\frac{\lambda}{3H(M)}, \frac{\epsilon}{2H(M)}, \delta/2)$
3: $\hat{F}_2 \leftarrow \text{AMS}(\epsilon, \delta/2)$
4: $S \leftarrow \{i \in \hat{S} : |g(\hat{v}_i) - g(\hat{v}_i + y)| \leq \epsilon g(\hat{v}_i + y)\}$,
5: for all $-\frac{\epsilon}{2H(M)}\sqrt{\hat{F}_2} \leq y \leq \frac{\epsilon}{2H(M)}\sqrt{\hat{F}_2}$
6: **return** (j, \hat{v}_j) , for all $j \in S$

LEMMA 23. *If an \mathbb{S} -normal function $g \in \mathcal{G}$ is not slow-dropping, then g is not 1-pass tractable.*

PROOF. Suppose g is not slow-dropping. Then there exist $0 < \alpha \leq 1$ and integer sequences $y_1, y_2, \dots \in \mathbb{N}$ and $x_1, x_2, \dots \in \mathbb{N}$, with y_k increasing, such that $x_k < y_k$ and $g(x_k) \geq y_k^\alpha g(y_k)$, for all $k \geq 1$. Furthermore, since g is \mathbb{S} -normal we may choose the sequences so that there exists a sub-polynomial function h satisfying

$$|g(x_k + y_k) - g(x_k)| > \frac{1}{h(y_k)} \min\{g(x_k), g(x_k + y_k)\}, \quad (1)$$

for all k . We claim that one can take $|g(x_k + y_k) - g(x_k)| > g(x_k)/h(y_k)$. If $g(x_k + y_k) < \frac{1}{2}g(x_k)$ then this is true because the constant function 2 is sub-polynomial. On the other hand, if $g(x_k + y_k) \geq \frac{1}{2}g(x_k)$ then replacing h by $2h$ in (1) does the job. Also, note that $1/h(x)$ is still a sub-polynomial function.

Let \mathcal{A} be a 1-pass (g, ϵ) -SUM algorithm with $\epsilon = (3h(n))^{-1}$. We will show that \mathcal{A} uses $\Omega(n^\alpha)$ bits on a sequence of g -SUM problems with $n = y_1, y_2, \dots$, hence g is not tractable. Let $n = y_k$ and $x = x_k$ and consider the following protocol for INDEX(n^α) using \mathcal{A} . Alice receives a set $A \subseteq [n^\alpha]$ and Bob receives an index $b \in [n^\alpha]$. Alice and Bob jointly create a notional stream D and run \mathcal{A} on it. Alice contributes n copies of i for each $i \in A$ to the stream and Bob contributes x copies of his index b . If $b \notin A$ there are $|A|$ items in the stream with frequency n and one with frequency x ; whereas if $b \in A$ then $|A| - 1$ frequencies are equal to n and one is $n + x$. Alice runs \mathcal{A} on her portion of the stream and sends the memory to Bob along with the value $|A|$. Bob completes the computation with his portion of the stream. Let \hat{G} be the value returned to Bob by \mathcal{A} . Bob decides that there is an intersection if $\frac{\hat{G}}{|A|g(n)+g(x)} \notin [1 - \epsilon, 1 + \epsilon]$.

With probability at least $2/3$, \hat{G} is a $(1 \pm \epsilon)$ -approximation to $g(V(D))$; suppose that this is so. If $b \in A$ then $g(V(D)) = (|A| - 1)g(n) + g(x + n)$ and otherwise the result is $g(V(D)) = |A|g(n) + g(x)$. Without loss of generality $g(n) < \epsilon g(x)$, thus the difference between these values is

$$\begin{aligned} |g(n) + g(x) - g(x + n)| &> |g(x) - g(x + n)| - \epsilon g(x) \\ &\geq 2\epsilon g(x) \\ &\geq \epsilon(g(x) + n^\alpha g(n)) \\ &\geq \epsilon(g(x) + |A|g(n)). \end{aligned}$$

Thus Bob's decision is correct and g -SUM inherits the $\Omega(n^\alpha)$ lower bound from INDEX(n^α). \square

LEMMA 24. *If an \mathbb{S} -normal function $g \in \mathcal{G}$ is not slow-jumping, then g is not 1-pass tractable.*

PROOF. Since g is not slow-jumping, there exist $\alpha > 0$, a strictly increasing sequence $y_1, y_2, \dots \in \mathbb{N}$, and a sequence

$x_1, x_2, \dots \in \mathbb{N}$ such that $x_k \leq y_k$ and $g(y_k) > s_k^{2+\alpha} x_k^\alpha g(x_k)$, where $s_k = \lfloor y_k/x_k \rfloor$. Without loss of generality y_k is strictly increasing. According to Lemma 23, we can assume that g is slow-dropping, since otherwise it is not 1-pass tractable. Hence, there exists $N \in \mathbb{N}$ such that for all $x \geq N$ and $r \leq x$ we have $\frac{1}{2}x^\alpha g(x) \geq g(r)$. If $x_k \leq N$ for all but finitely many k , we may assume this holds for all k . Otherwise, by taking a subsequence, we can assume $x_k > N$ holds for all k .

Let $r_k = y_k - s_k x_k$ be the sequence remainders. Before proceeding with the reduction we will establish the bound $2g(r_k) \leq g(y_k)$. If $x_k \leq N$, for all k , then s_k is unbounded while $x_k g(x_k)$ is bounded away from zero. This means that $g(y_k)$ is unbounded while $g(r_k)$ is bounded and we can assume y_1 is large enough that $g(y_k) \geq 2g(r_k)$ holds for all k . Otherwise $x_k > N$, for every k , hence $2g(r_k) \leq x_k^\alpha g(x_k) \leq g(y_k)$, for each k , from the definition of N .

Let \mathcal{A} be a (g, ϵ) -SUM algorithm with accuracy $\epsilon = 1/12$. Consider the DISJ-IND(n, t) problem, where $t = s_k$ and $n = s_k^{2+\alpha} x_k^\alpha$. Denote $x := x_k, y := y_k, s := s_k$ and $r := r_k$. The $t + 1$ players receive sets $A_1, A_2, \dots, A_t \subseteq [n]$ and an index $b \in [n]$. As before, the players jointly run \mathcal{A} on a notional stream and share $|A_i|$. Each player i places x copies of each $j \in A_i$ into the stream except for the final player who places r copies of his index b . Let $n' = \sum_i |A_i|$ be the total size of the t players' sets. On an intersecting instance the result of g -SUM is $a_1 := (n' - t)g(x) + g(y)$, and on a disjoint instance the result is $a_2 := n'g(x) + g(r)$.

The difference is

$$\begin{aligned} a_1 - a_2 &\geq g(y) - g(r) - tg(x) > \frac{1}{2}g(y) - tg(x) \\ &\geq \frac{1}{6}(g(y) + (2s^{2+\alpha} x^\alpha)g(x)) - sg(x) \\ &= \frac{1}{6}(g(y) + (s^{2+\alpha} x^\alpha + s^{2+\alpha} x^\alpha - 6s)g(x)) \\ &\geq 2\epsilon(g(y) + s^{2+\alpha} x^\alpha g(x)) \\ &= 2\epsilon(n'g(x) + g(y)), \end{aligned} \quad (2)$$

where the last inequality holds for all sufficiently large n . Thus, the index player can correctly distinguish which case has occurred, and the algorithm requires $\Omega(n/t^2) = \Omega(y^\alpha)$ bits. \square

LEMMA 25. *If an \mathbb{S} -normal function $g \in \mathcal{G}$ is not predictable, then g is not 1-pass tractable.*

PROOF. Since g is not predictable, there exists a sub-polynomial function ϵ and constant $\gamma > 0$ such that some infinite sequence x_k, y_k satisfies the following:

- $x_k \rightarrow \infty$ and $y_k \in [1, x_k^{1-\gamma})$,
- $x_k + y_k \notin \delta_{\epsilon(x_k)}(g, x_k)$, and
- $x_k^\gamma g(y_k) < g(x_k)$.

The proof is by a reduction from INDEX(n) with $n = \epsilon(x_k)x_k^\gamma/4$. Suppose \mathcal{A} is an algorithm for $(g, \epsilon/4)$ -SUM. Alice receives a set $A \subseteq [n]$ and Bob receives an index $b \in [n]$. Alice adds y_k copies of i to the stream, for each $i \in A$, runs \mathcal{A} on her portion of the stream, and sends the contents of the memory to Bob. Bob adds x_k copies of b to the stream and completes the computation.

The stream has $|A|$ frequencies equal to y_k and one equal to x_k , if there is no intersection, or $|A| - 1$ equal to y_k and

one equal to $x_k + y_k$, if there is an intersection. Recall that $x_k + y_k \notin \delta_\epsilon(g, x_k)$, hence $|g(x_k) - g(x_k + y_k)| > \epsilon g(x_k)$, and by construction

$$|A|g(y_k) \leq \frac{\epsilon(x_k)}{4} x_k^\gamma g(y_k) \leq \frac{\epsilon(x_k)}{4} g(x_k).$$

Therefore, Bob can correctly distinguish whether $b \in A$ when \mathcal{A} yields a $(1 \pm \epsilon/4)$ -approximation. Thus $(g, \epsilon/4)$ -SUM requires $\Omega(n)$ bits. \square

4.5 Two Pass Lower Bounds

THEOREM 26. *If a \mathbb{P} -normal function $g \in \mathcal{G}$ is tractable then g is slow-dropping and slow-jumping.*

PROOF. Immediate from Lemmas 27 and 28. \square

LEMMA 27. *If a \mathbb{P} -normal function $g \in \mathcal{G}$ is not slow-dropping, then g is not $O(1)$ -pass tractable.*

PROOF. Suppose \mathcal{A} is a $p = O(1)$ pass algorithm that solves (g, ϵ) -SUM and g does not satisfy slow-dropping condition. Then there exists $\alpha > 0$ such that for any $N > 0$, there exists $y > N$ and $x < y$ satisfying $g(y) < g(x)/y^\alpha$. Separately, since g is \mathbb{P} -normal, there exists $\beta > 0$ such that for any $N > 0$ if g has an α -period $y > N$ then there is an α -period $y > N$ and $x < y$ such that $g(y) \leq g(x)/y^\alpha$ and $|g(x+y) - g(x)| > y^\beta \min(g(x), g(x+y))$.

Let $\gamma = \min(\alpha, \beta)$. Consider the following protocol for DISJ($n, 2$), where $n = y^{\gamma/2}$.

First, consider $g(x+y) \leq g(x)$. Since $|g(x+y) - g(x)| > y^\gamma g(x+y)$ and $y^\gamma g(x+y) \geq g(x+y)$, then $g(x) \geq y^\gamma g(x+y)$. The players jointly create a stream where Player 1 inserts x copies of each element of her set S_1 into the stream, and Player 2 inserts y copies of every element a not in her set S_2 . First Player 1 creates her portion of the stream, runs the first pass of \mathcal{A} on it, and sends the memory to Player 2. She completes the first pass with her portion of the stream, and returns the memory to Player 1. The players continue in this way for a total of p passes over the stream.

Let V denote the frequency vector of the resulting stream. If there is an intersection, let $S_1 \cap S_2 = \{a\}$. Then $S_1 \cap S_2 = S_1 \setminus \{a\}$, and $g(V)$ is $r_1 = (|S_1| - 1)g(x+y) + (n - |S_2| - |S_1|)g(y) + g(x) + g(y)$. If there is no intersection, the value of $g(V)$ is $r_2 = (|S_1| - 1)g(x+y) + (n - |S_2| - |S_1|)g(y) + g(x+y)$. Notice that

$$\frac{|r_2 - r_1|}{r_1} \geq \frac{|g(x+y) - g(x)|}{g(x)} \geq \frac{1}{2},$$

for sufficiently large n . For any $\epsilon < 1/2$, \mathcal{A} is able to distinguish the 2 cases, which gives a lower bound on the memory bits used by \mathcal{A} is $\Omega(n/p)$.

The case $g(x+y) > g(x)$ is the same except Player 2 inserts y copies of each element that is in her set S_2 . Thus, \mathcal{A} uses $\Omega(n/p)$ bits of memory, hence g is not $O(1)$ -pass tractable. \square

LEMMA 28. *If a \mathbb{P} -normal function $g \in \mathcal{G}$ is not slow-jumping, then g is not $O(1)$ -pass tractable.*

PROOF. Suppose \mathcal{A} is a $p = O(1)$ pass algorithm for (g, ϵ) -SUM. If g is not slow-jumping, then there exists $\alpha > 0$ such that for any $N > 0$, there exists $x < y \in \mathbb{N}$ and $y \geq N$ but $g(y) > \lfloor y/x \rfloor^{2+\alpha} x^\alpha g(x)$. Notice that for $y > x$ we have $\lfloor y/x \rfloor \geq y/2x$, so by adjusting α we can assume $g(y) > (y/x)^2 y^\alpha g(x)$.

We can further assume g is slow dropping, since otherwise g is not $O(1)$ -pass tractable by Lemma 27. Thus, there exists a nondecreasing sub-polynomial function $h(x) > 1$ such that $g(x) \leq h(y)g(y)$.

Consider an instance $A_1, A_2, \dots, A_t \subseteq [n]$ of DISJ(n, t), where $t = \lceil y/x \rceil$ and $n = (\frac{y}{x})^2 \frac{y^\alpha}{2h(y)}$. Each of the first $t-1$ players, inserts x copies of her elements into the stream and the t th player inserts $y - (t-1)x < x$ copies of her elements into the stream. The players pass the memory and repeat to run the p passes of \mathcal{A} as before. Notice that $g(y - (t-1)x) \leq h(x)g(x)$ by the slow dropping condition.

Let $n' = \sum_i |A_i|$, and let V be the frequency vector of the stream created. If there is no intersection, the value of $g(V)$ is

$$\begin{aligned} (n' - |A_t|)g(x) + |A_t|g(y - (t-1)x) &\leq n'h(x)g(x) \\ &\leq nh(x)g(x) \leq \frac{1}{2}g(y). \end{aligned}$$

On the other hand, if there is an intersection then y is the frequency of some item in the stream, and therefore $g(V) \geq g(y)$. Thus \mathcal{A} distinguishes between the cases. The algorithm uses $\Omega(n/pt^2) = \Omega(y^{\alpha/2}/p)$ bits of memory, which proves that g is not p -pass tractable. \square

4.6 Zero-One Law Proofs

THEOREM 2 (1-PASS ZERO-ONE LAW). *A function $g \in \mathcal{G}$ is 1-pass tractable and \mathbb{S} -normal if and only if it is slow-jumping, slow-dropping, and predictable.*

PROOF. The lower bound is proved as Theorem 26, hence the algorithm remains. With the Recursive Sketch of Theorem 13, it is enough to show that Algorithm 2 is a $(g, \lambda, \epsilon, \delta)$ -heavy hitters algorithm for $\lambda = \epsilon^2 / \log^3 n$ and $\delta = 1 / \log n$.

By Lemma 18, the Count Sketch used by Algorithm 2 returns a list of pairs (i_j, \hat{v}_{i_j}) containing all of the $\lambda/3H(M)$ -heavy elements for F_2 . By definition, any item i_j that survives the pruning stage has $|g(\hat{v}_{i_j}) - g(v_{i_j})| \leq \epsilon g(v_{i_j})$. Hence, it only remains to show that every (g, λ) -heavy hitter survives the pruning stage.

Suppose i' is the index of the (g, λ) -heavy hitter that minimizes $r_{\epsilon/2}(v_{i'})$. Now suppose that we insert a new item i'' in the stream with frequency $v_{i''} = r_{\epsilon/2}(v_{i'}) + 1$. Then

$$\frac{g(v_{i'})}{H(M)} \leq g(v_{i''}) \leq H(M)g(v_{i'}),$$

where the first inequality follows from Lemma 21 and the second because g is slow-dropping. Thus, this item is a $(g, \frac{\lambda}{3H(M)})$ -heavy hitter and the constant term on the parameter b for the Count Sketch can be chosen to guarantee additive error no more than $v_{i''}/3 \leq r_{\epsilon/2}(v_{i'})/3$. The same guarantee holds for the stream without i'' , thus for every (g, λ) -heavy hitter i_j we have

$$|\hat{v}_{i_j} - v_{i_j}| \leq \frac{1}{3}r_{\epsilon/2}(v_{i'}) \leq \frac{1}{3}r_{\epsilon/2}(v_{i_j}),$$

where the last inequality holds by the choice of i' . Furthermore, for every $-r_{\epsilon/2}(v_{i'})/3 \leq y \leq r_{\epsilon/2}(v_{i'})/3$ we have $|(\hat{v}_{i_j} + y) - v_{i_j}| \leq r_{\epsilon/2}(v_{i_j})$ and thus

$$\begin{aligned} |g(\hat{v}_{i_j}) - g(\hat{v}_{i_j} + y)| &\leq |g(\hat{v}_{i_j}) - g(v_{i_j})| \\ &\quad + |g(v_{i_j}) - g(\hat{v}_{i_j} + y)| \\ &\leq \epsilon g(\hat{v}_{i_j} + y). \end{aligned}$$

With probability at least $1 - 2\delta/2 = 1 - \delta$ both the Count Sketch and the AMS approximation of F_2 meet their obligations, in this case the output is correct. By the guarantee of CountSketch, we can safely assume that $r_{\epsilon/2}(v'_i) \geq \epsilon/2H(M)\sqrt{\widehat{F}_2}$, heavy hitters will survive the pruning stage. \square

THEOREM 3 (2-PASS ZERO-ONE LAW). *A function $g \in \mathcal{G}$ is 2-pass tractable and \mathbb{P} -normal if and only if it is slow-dropping and slow-jumping. Furthermore, every slow-dropping and slow-jumping \mathbb{S} -normal function is also 2-pass tractable.*

PROOF. The lower bound is proved as Theorem 26. The upper bound is governed by Proposition 10 and Theorem 19. \square

Here are a few examples. The functions $x^2 \lg(1+x)$, $(2 + \sin \log(1+x))x^2$, and $e^{\log^{1/2}(1+x)}$ are all 1-pass tractable because they are all slow-dropping, slow-jumping, and predictable. On the other hand, $1/x$ is not slow-dropping, x^3 is not slow-jumping, and $(2 + \sin \sqrt{x})x^2$ is not predictable, so none of these functions is 1-pass tractable. The last of the three is, however, slow-jumping and slow-dropping, and hence it is 2-pass tractable.

5. NEARLY PERIODIC FUNCTIONS

Nearly periodic functions are highly constrained to *almost* repeat themselves like a periodic function. These functions admit large changes in value that would imply a large lower bound on their space complexities, if they did not satisfy the many constraints.

5.1 Example Nearly Periodic Function

Constructing a tractable nearly periodic function turns out to be a non-trivial exercise. This section provides such a construction.

DEFINITION 29. *Let $x = \sum_{j=0}^{\infty} a_j 2^j \in \mathbb{N}$, for $a_j \in \{0, 1\}$, be the binary expansion of the integer x . Define $i_x := \min\{j : a_j = 1\}$, the location of the first non-zero bit of x . The function is $g_{np}(x) = 2^{-i_x}$, for $x > 0$, and $g_{np}(0) = 0$.*

PROPOSITION 30. *g_{np} is \mathbb{S} -nearly periodic.*

PROOF. Since there is an infinite sequence $\{1, 2, 4, \dots, 2^n \dots\}$ such that $g_{np}(x) = 1/x$, the first condition of near-periodicity is satisfied. It remains to show that the second condition is also satisfied.

Let $\gamma > 0$ be a constant. Consider any integer $x > 0$. By construction $g_{np}(x) = 2^{-i_x}$. Let $y > x$ be another integer. Similarly, we have $g_{np}(y) = 2^{-i_y}$. If $g_{np}(x)/g_{np}(y) = 2^{i_y - i_x} \geq y^\gamma$, then $i_y - i_x \geq \gamma \lg y$. Choose $N = \lceil 2^{1/\gamma} \rceil$. For any $y > N$, we have $i_y - i_x > 1$ and therefore $i_{x+y} = i_x$. Thus, $g_{np}(x+y) = g_{np}(x)$. \square

PROPOSITION 31. *g_{np} is 1-pass tractable.*

PROOF. We demonstrate that one can find (g_{np}, λ) -heavy hitters in $\text{poly}(\lambda^{-1} \log n \log M)$ space. It is sufficient to demonstrate an algorithm that will find a single (g_{np}, λ) -heavy hitter, since we can reduce the heavy hitters problem to this case by hashing the stream $O(\lambda^{-2})$ ways and running this algorithm on each substream.

Suppose that j^* is the single (g_{np}, λ) -heavy hitter with frequency x . Let $v_1, v_2, \dots, v_n \geq 0$ be the frequencies in the stream, and let $U = \{j : i_{v_j} \leq i_x\}$. By definition $g_{np}(v_j) \geq g_{np}(x)$ for $j \in U$, hence $|U| \leq 1 + \lambda^{-1} \leq 2\lambda^{-1}$.

Let $C = O(\lambda^{-2})$. Apply a uniform hash function $h : [n] \rightarrow [C]$ to separate the stream into C substreams S_1, S_2, \dots, S_C . With constant probability, no two elements of U are in the same substream, so suppose that this happens.

On each substream S_k with frequencies $v_1^{(k)}, \dots, v_n^{(k)}$ we run the following algorithm $D = O(\log n)$ times independently in parallel: (i) Sample pairwise independent Bernoulli $(1/2)$ random variables X_1, \dots, X_n . (ii) Compute $m = \sum_j X_j v_j^{(k)}$ and i_m . (iii) Output 2^{-i_m} .

Consider the set of D values output by this algorithm. If there is a single item j^* with minimum $i^* := i_{v_{j^*}^{(k)}}$, then a Chernoff bound implies that very nearly $D/2$ values are equal to 2^{-i^*} , and this is the maximum value among all of the pairs. In this case, the label j^* can be found in post-processing by binary search. Let $X_{j,\ell}$ denote the value of the j th Bernoulli variable on the ℓ th trial. Specifically, one finds the set $M \subseteq [D]$ of trials for which 2^{-i_m} is equal to the maximum among the D values, and with high probability, only j^* will satisfy $X_{j,\ell} = 1$ for all $\ell \in M$ and $X_{j,\ell} = 0$ for all $\ell \in [D] \setminus M$. We can detect the case where there is more than one item with minimum $i_{v_j^{(k)}}$ because either the number of maximizing values will be too large or the binary search will fail to yield a unique element.

The single-heavy-hitter algorithm now outputs the pair $(j^*, 2^{-i^*})$ if the number of maximizing values is correct and the binary search yields a unique element or nothing if either of those conditions fails. With the extra hashing step mentioned at the beginning this yields a 1-pass $O(\lambda^{-4} \log n \log M)$ -space (g_{np}, λ) -heavy hitters algorithm, thus g_{np} is 1-pass tractable. \square

5.2 More Lower Bounds

The set of nearly-periodic function defeat the standard reduction from the DISJ and INDEX problems. Theorem 35 shows how one can use the SHORTLINEARCOMBINATION problem, of Definition 14, to provide space lower bounds on g -SUM algorithms for some nearly periodic functions g .

DEFINITION 32. *For any $N > 0$ and non-increasing sub-polynomial function $h(x)$, define the (N, h) -dropping set of f to be $\mathcal{D}_{N,h}(f) = \{x \mid 1 \leq x \leq N, f(x) \leq h(N)/N\}$.*

PROPOSITION 33. *Let f be a nearly-periodic function, then for all $n_0 > 0$, there exists $N > n_0$ and a sub-polynomial function h such that $|\mathcal{D}_{N,h}| > 0$.*

PROOF. Choose $N = n_0 + 1$, $h(x) = f(1)N$, i.e. $h(x)$ is a constant function. Then 1 in $\mathcal{D}_{N,h}$ because $1 \leq N$ and $f(1) = h(N)/N$. \square

DEFINITION 34. *An α -indistinguishable frequency set of $[n]$ is a tuple (s, d) , where $s \subseteq [n]$ is a subset and $d \in [n]$, $d \notin s$ is a integer such that (s, d) -DIST problem requires $\Omega(n^\alpha)$ space.*

THEOREM 35. *Let function $f : \mathbb{R} \rightarrow \mathbb{R}^+$, symmetric, nearly-periodic and $f(0) = 0$. If there exists a non-increasing sub-polynomial function $h(x)$ and constants $\alpha > 0$, $0 < \delta < 1$, such that for any $n_0 > 0$, there exists $N > n_0$ and a integer set $|S| > 0$ satisfying the follows,*

1. $0 < |\mathcal{D}_{N,h}(f)| \leq N^{1-\delta}$;
2. $S \subseteq |\mathcal{D}_{N,h}(f)|$ and $d \in [N]$ such that (S, d) is a α -indistinguishable frequency set of N ;

$$3. f(d) \geq h(d),$$

then for any $n_0 > 0$, there exists a stream of domain $N > n_0$, any one-pass randomized streaming algorithm \mathcal{A} that outputs a $1 \pm \epsilon(N)$ approximation of f -SUM with probability at least $2/3$, requires space $\Omega(n_0^\alpha)$ bits, where $\epsilon(x) < 1$ is a sub-polynomial function.

PROOF. Suppose we have an algorithm \mathcal{A} that gives a $1 \pm \epsilon(N)$ approximation to f -SUM with probability at least $2/3$. We can now use \mathcal{A} to construct a protocol that solves (S, d) -DIST problem on domain N . The algorithm is straightforward: run \mathcal{A} on the input stream. If the result \mathcal{A} outputs is $\leq h(d)/N^\delta$, then there is no frequency of absolute value d in the stream. This is so since the g -SUM is at most $\sum_{a \in S} f(a) \leq N^{1-\delta}/N$. Otherwise, if the output is at least $h(d)(1 - \epsilon(N))$. By the guarantee of \mathcal{A} , it can distinguish the 2 cases. The space needed of \mathcal{A} inherits the lower bound of (S, d) -DIST problem, thus $\Omega(N^\alpha)$ bits. \square

6. REFERENCES

- [1] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Proceedings of the 28th annual ACM Symposium on Theory of Computing*, pages 20–29. ACM, 1996.
- [2] Alexandr Andoni, Robert Krauthgamer, and Ilya Razenshteyn. Sketching and embedding are equivalent for norms. In *Proceedings of the 47th annual ACM Symposium on Theory of Computing*. ACM, 2015.
- [3] Ziv Bar-Yossef. *The complexity of massive data set computations*. PhD thesis, University of California at Berkeley, 2002.
- [4] Ziv Bar-Yossef, TS Jayram, Ravi Kumar, and D Sivakumar. An information statistics approach to data stream and communication complexity. In *Proceedings of the 43rd annual IEEE Symposium on Foundations of Computer Science*, pages 209–218. IEEE, 2002.
- [5] Vladimir Braverman and Stephen R Chestnut. Universal sketches for the frequency negative moments and other decreasing streaming sums. *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, page 591, 2015.
- [6] Vladimir Braverman and Rafail Ostrovsky. Zero-one frequency laws. In *Proceedings of the 42nd annual ACM Symposium on Theory of Computing*, pages 281–290. ACM, 2010.
- [7] Vladimir Braverman and Rafail Ostrovsky. Generalizing the layering method of Indyk and Woodruff: Recursive sketches for frequency-based vectors on streams. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 58–70. Springer, 2013.
- [8] Amit Chakrabarti, Subhash Khot, and Xiaodong Sun. Near-optimal lower bounds on the multi-party communication complexity of set disjointness. In *Proceedings of the 18th annual Conference on Computational Complexity*, pages 107–117. IEEE, 2003.
- [9] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Automata, Languages and Programming*, pages 693–703. Springer, 2002.
- [10] Stephen R Chestnut. *Stream sketches, sampling, and sabotage*. PhD thesis, Johns Hopkins University, 2015. <https://jscholarship.library.jhu.edu/handle/1774.2/37935>.
- [11] André Gronemeier. Asymptotically optimal lower bounds on the NIH-multi-party information complexity of the AND-function and disjointness. In *STACS 2009: 26th International Symposium on Theoretical Aspects of Computer Science*, volume 3 of *LIPICs. Leibniz Int. Proc. Inform.*, pages 505–516. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, 2009.
- [12] Sudipto Guha and Piotr Indyk. Problem 5: Characterizing Sketchable Distances. <http://sublinear.info/5>, 2006.
- [13] Sudipto Guha, Piotr Indyk, and Andrew McGregor. Sketching information divergences. In *Learning Theory*, pages 424–438. Springer, 2007.
- [14] Piotr Indyk. Stable distributions, pseudorandom generators, embeddings and data stream computation. In *Proceedings of the 41st annual IEEE Symposium on Foundations of Computer Science*, pages 189–197. IEEE, 2000.
- [15] Piotr Indyk and David Woodruff. Optimal approximations of the frequency moments of data streams. In *Proceedings of the 37th annual ACM Symposium on Theory of Computing*, pages 202–208. ACM, 2005.
- [16] Eaman Jahani, Michael J Cafarella, and Christopher Ré. Automatic optimization for mapreduce programs. *Proceedings of the VLDB Endowment*, 4(6):385–396, 2011.
- [17] T. S. Jayram. Hellinger strikes back: A note on the multi-party information complexity of AND. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 562–573, 2009.
- [18] Ilan Kremer, Noam Nisan, and Dana Ron. On randomized one-round communication complexity. *Computational Complexity*, 8(1):21–49, 1999.
- [19] Kyong-Ha Lee, Yoon-Joon Lee, Hyunsik Choi, Yon Dohn Chung, and Bongki Moon. Parallel data processing with mapreduce: a survey. *ACM SIGMOD Record*, 40(4):11–20, 2012.
- [20] Yi Li, Huy L Nguyen, and David P Woodruff. Turnstile streaming algorithms might as well be linear sketches. In *Proceedings of the 46th annual ACM Symposium on Theory of Computing*, pages 174–183. ACM, 2014.
- [21] Yi Li and David P Woodruff. A tight lower bound for high frequency moment estimation with small error. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 623–638. Springer, 2013.
- [22] Sai Wu, Feng Li, Sharad Mehrotra, and Beng Chin Ooi. Query optimization for massively parallel data processing. In *Proceedings of the 2nd ACM Symposium on Cloud Computing*, page 12. ACM, 2011.

APPENDIX

A. THE CASE OF $G(0) \neq 0$

When $g(0) = 0$, we have to re-address the lower bound of using INDEX since the elements not contributing to the frequency vector still contribute to the g -SUM. Therefore, we re-define the nearly periodic functions in the next section. We further show that functions crossing the axis are not 1-pass tractable and functions with zero points are not tractable unless they are periodic in Section A.2. We prove the same 1-pass zero one law in the remaining 2 sections. The proof turns out to be very similar to the $g(0) = 0$ cases with small number of additional tweaks. Note that we only provide the turnstile model lower bounds.

A.1 Redefinition of Nearly Periodic

For the $g(0) \neq 0$ case, we have the following definition of nearly-periodic.

DEFINITION 36. *Given a set of functions \mathcal{S} , call $g(x)$ \mathcal{S} -nearly periodic, if the following two conditions are satisfied.*

1. *There exists $\alpha > 0$ such that for any constant $N > 0$ there exists $x, y \in \mathbb{N}$, $x < y$ and $y \geq N$ such that $g(y) \leq g(x)/y^\alpha$. Call such a y α -period of g ;*
2. *For any $\alpha > 0$ and any error function $h \in \mathcal{S}$ there exists $N_1 > 0$ such that for all α -periods $y \geq N_1$ all $x < y$ such that $g(y)y^\alpha \leq g(x)$ we have $|g(x) - g(x - 2y)| \leq \min\{g(x), g(x - 2y)\}h(y)$.*

A function g is \mathcal{S} -normal if it is not \mathcal{S} -nearly periodic.

Let all functions be $\mathcal{G}_0^* = \{g : \mathbb{N} \rightarrow \mathbb{R}^+, g(x) = g(-x), \text{ and } g(0) = 1\}$.

A.2 Crossing the Axis

If $g \in \mathcal{G}_0^*$ and there exist $x, y \in \mathbb{N}$ such that $g(x) > 0 > g(y)$, we have the following lemma.

LEMMA 37. *Let $g \in \mathcal{G}_0^*$, if $g(1) < 0$, then any algorithm solves g -SUM requires $\Omega(n)$ space.*

PROOF. Consider the following INDEX(n) reduction. Alice and Bob jointly create a stream as following. Let A be Alice's set. $n_1 = |A|$. Let $n_2 = \lfloor \frac{n - n_1}{-g(1)} \rfloor$. Let $C = (n_1 + n_2)g(1) + (n - n_1)$. We have $0 \leq C < g(1)$. Alice choose a domain $[n + n_2]$ for the algorithm and output 1 copy of each of her element. Bob receives the memory content of the algorithm from Alice as well as n, n_1, n_2 . Bob outputs -1 copy of his index and also adds to the stream one copy of each of the following elements: $n + 1, n + 2, \dots, n + n_2$.

If there is an intersection, the result of g -SUM is

$$(n - n_1) + (n_1 + n_2 - 1)g(1) + g(0) = C + g(0) - g(1);$$

If there is no intersection, the result is

$$(n - n_1 - 1) + (n_1 + n_2)g(1) + g(1) = C + g(1) - g(0).$$

Then a constant approximation algorithm can distinguish the two cases. The algorithm inherits an $\Omega(n)$ bound from INDEX. \square

LEMMA 38. *Let $g \in \mathcal{G}_0^*$ and $g(1) > 0$, if there exists $y \in \mathbb{N}$ such that $g(y) < 0 < g(1)$, then there exists $z \in \mathbb{N}$ such that $g(1 + z) \neq g(1 - z)$.*

PROOF. If for every z , $g(1 + z) = g(1 - z)$, we have $g(2) = g(0) = 1$ and for any w , $g(w + 2) = g(w)$. Therefore 2 is a period of the function. Now, consider the following cases, 1) if y even, then $g(y) = g(2\lfloor y/2 \rfloor) = g(0)$ (2 is a period), which is a contradiction that $g(y) < 0$; 2) if y odd, then $g(y) = g(2\lfloor y/2 \rfloor + 1) = g(1) > 0$, contradicting $g(y) < 0$. \square

PROPOSITION 39. *Let $g \in \mathcal{G}_0^*$ and $g(1) > 0$, if there exists $y \in \mathbb{N}$ such that $g(y) < 0$, then any algorithm solves g -SUM requires $\Omega(n)$ space.*

PROOF. Consider the following reduction from INDEX(n). Let $n' = \lfloor -(n - 1)g(1)/g(y) \rfloor$ and $z \in \mathbb{N}$ be the integer given by Lemma 38. Let $C = (n - 1)g(1) + n'g(y)$. We have $0 < C < -g(y)$. Alice and Bob jointly create a data stream in domain $n + n'$: Alice outputs 1 copy of every element in her set, and -1 copy of every element not in her set. Bob outputs y copy of each of the elements $i + 1, i + 2, \dots, i + n'$ and z copy of his index. If there is no intersection, the result of g -SUM is

$$(n - 1)g(x) + n'g(y) + g(z - 1) = C + g(z - 1).$$

If there is an intersection, the result of g -SUM is

$$(n - 1)g(x) + n'g(y) + g(z + 1) = C + g(z + 1)$$

Therefore any constant approximation algorithm can distinguish the two cases, implying an $\Omega(n)$ bound for the memory of the algorithm. \square

PROPOSITION 40. *Let $g \in \mathcal{G}_0^*$ if $g(x) = 0$ and $g(2x) \neq g(0)$ then g is not 1-pass tractable.*

PROOF. Consider the INDEX(n) reduction. Alice and Bob jointly create a stream as following, Alice output x for every element in her set, $-x$ for every element not in her set. Bob output x for his index. If there is an intersection, the g -SUM result is $g(2x)$. If there is no intersection, the result is $g(0)$. \square

PROPOSITION 41. *If function $g \in \mathcal{G}_0^*$ with $g(x) = 0$ for some $x \in \mathbb{N}$ is tractable, then g is periodic.*

PROOF. By reduction from INDEX(n), Alice output x for each of her elements and $-x$ for each of her non-elements. Bob output $k + x$ for his index. If there is an intersection, the g -SUM result is $g(k + 2x)$. If there is no intersection, g -SUM is $g(k)$. If $g(k + 2x) \neq g(k)$ for any k , a constant approximation algorithm can distinguish the 2 cases. Therefore, $2x$ is a period of the function. \square

Therefore, we only consider the following set of functions,

$$\mathcal{G}_0 = \{g : \mathbb{N} \rightarrow \mathbb{R}^+, g(x) = g(-x) > 0, \text{ and } g(0) = 1\}$$

A.3 Lower Bounds

THEOREM 42. *If an \mathcal{S} -normal function $g \in \mathcal{G}_0$ does not satisfy slow dropping, then g is not 1-pass tractable.*

PROOF. The reduction to INDEX(n^α) is similar to the case of $g(0) = 0$. Now Alice outputs n copies of the elements in A and also outputs $-n$ copies of i if $i \notin A$. Bob output $x - n$ copies of his index. If there is an intersection the g -SUM result is $g(V(D)) = (n^\alpha - 1)g(n) + g(x)$ and otherwise is $g(V(D)) = (n^\alpha - 1)g(n) + g(x - 2n)$. Since g is not nearly periodic, the algorithm is able to distinguish the two cases. \square

THEOREM 43. *If an \mathbb{S} -normal function $g \in \mathcal{G}_0$ is not slow-jumping, then g is not 1-pass tractable.*

PROOF. The reduction to DISJ+IND is the same to the case of $g(0) = 0$.

Now by slow dropping, there exists a sub-polynomial function $h(x)$ such that $g(0) \leq g(x)h(x)$. The g -SUM result differs from $g(0) = 0$ case by $(n - n')g(0)$ (or $(n - n' + t)g(0)$), which is comparable to $(n - n')g(x)$. Therefore, we can apply the same analysis. \square

THEOREM 44. *If an \mathbb{S} -normal function $g \in \mathcal{G}_0$ is not predictable, then g is not 1-pass tractable.*

PROOF. Still use the same reduction to INDEX with the $g(0) = 0$ case. Now the g -SUM results differs by $(n - |A|)g(0)$. By slow dropping, there exists a sub-polynomial function $h(x)$ such that $g(0) \leq h(y)g(y)$. Therefore, all the analyses remain the same. \square

A.4 Upper Bound

The algorithm should still work by the following 2 lemmas.

LEMMA 45. *Let $g \in \mathcal{G}_0$ be an \mathbb{S} -normal function that is slow-jumping and slow-dropping. There exists a sub-polynomial function h such that for any $D \in \mathcal{D}(n, m)$ with frequencies v_1, v_2, \dots, v_n , if $g(v_i) \geq \lambda \sum_j g(v_j)$ then $v_i^2 \geq \frac{\lambda}{h(|v_i|)} \sum_{|v_j| < |v_i|} v_j^2$.*

PROOF. The proof is the same as in the $g(0) = 0$ case with additional care of $g(0) = 1$. \square

LEMMA 46. *If $g \in \mathcal{G}_0$ is slow-dropping, slow-jumping, and predictable, then there is a sub-polynomial function h such that for all $y \in [r_\epsilon(x) + 1, x/h(x)]$ we have $g(y) \geq g(x)/h(x)$.*

PROOF. The proof is the same as in the $g(0) = 0$ case with additional care of $g(0) = 1$. \square

B. LOWER BOUND FOR DISJ+IND

THEOREM 47. *The randomized one-way communication complexity of DISJ+IND(n, t) is $\Omega(n/t \log n)$.*

PROOF. We give a reduction from DISJ($n, t + 1$). Let \mathcal{P} be any randomized protocol for DISJ+IND(n, t). Run in parallel $\ell = \lceil 96 \log n \rceil$ independent copies of \mathcal{P} through the first t players. This produces ℓ transcripts. Player $t + 1$ now takes the each of the transcripts and computes the final value of each once for every element he holds, as if it was the only element he held. No communication is need for this part because \mathcal{P} is a one-way protocol and $t + 1$ is the final player.

Player $t + 1$ then takes the majority vote among the independent copies of \mathcal{P} for each element. If any vote signals an intersection then he reports intersection; otherwise he reports disjoint.

If every one of the $|A_{t+1}| \leq n$ majorities is correct then the final player's report is correct. Let X_i , for $i \in A_{t+1}$, be the number among the ℓ copies of \mathcal{P} with the correct outcome when the final player completes protocol using $i \in A_{t+1}$. Then X_i is Binomially distributed from ℓ trials with success probability at least $2/3$. Using a Chernoff Bound we find

$$\begin{aligned} P(X_i \leq \ell/2) &= P\left(X_i \leq (1 - \frac{1}{4})\frac{2}{3}\ell\right) \leq \exp\left\{\frac{-1}{32}\mu\right\} \\ &\leq \exp\left\{\frac{-1}{32} \cdot \frac{2}{3}\ell\right\} \leq \frac{1}{n^2}. \end{aligned}$$

Thus, with probability at least $1 - \frac{1}{n}$ every majority vote is correct, hence our DISJ(n, t) protocol is correct for $n \geq 3$.

Let T_1, T_2, \dots, T_ℓ be the transcripts. The total cost of this DISJ protocol is $\sum_{i=1}^\ell |T_i|$. Since DISJ($n, t + 1$) requires $\Omega(n/t)$ bits of communication, at least one of the protocols has length $\Omega(n/t\ell) = \Omega(n/t \log n)$, hence $|\mathcal{P}| = \Omega(n/t \log n)$ bits of communication. \square

C. LOWER BOUND FOR SHORTLINEARCOMBINATION

C.1 The 3-frequency Distinguishing Problem

We first define the ShortLinearCombination problem for three frequencies, which we call (a, b, c) -DIST for short.

DEFINITION 48. *A stream S with frequency vector v is given to a one-pass streaming algorithm. v is promised to be from $V_0 = \{-a, a, -b, b, 0\}^n$ or $V_1 = \{EMB(v, i, e) \mid v \in V_0, i \in [n], e \in \{-c, c\}\}$, where $a, b, c > 0$ and*

$$EMB(v, i, e)_j = \begin{cases} v_i & i \neq j \\ e & i = j \end{cases}$$

(in other words, V_1 is given by replacing a coordinate of a vector v from V_0 with $-c$ or c). The (a, b, c) -DIST problem is for the algorithm to distinguish whether $v \in V_0$ or $v \in V_1$.

To study this problem, we first consider the special case where $c = \gcd(a, b) = 1$.

PROPOSITION 49. *If $\gcd(a, b) = 1$, then any randomized algorithm solving $(a, b, 1)$ -DIST with probability at least $2/3$ requires $\Omega(n/\max(a, b)^2)$ bits of memory.*

PROOF. We use the same notation as in [4], and in particular refer to that work for background on information complexity. Without loss of generality, we assume $a > b$. Consider a communication problem in which Alice is given a vector v_1 and Bob is given a vector v_2 , for which $v = v_1 - v_2 \in V_0 \cup V_1$. The players are asked to solve the $(a, b, 1)$ -DIST on v . This problem is OR-decomposable with primitive DIST(x, y), where DIST(x, y) = 1 if and only if $|x - y| = 1$. Consider a randomized protocol Π which succeeds with probability at least $2/3$ in solving this problem, and suppose Π has the minimum communication cost of all such protocols. We now call $(a, b, 1)$ -DIST f for short. The following is well known (see, e.g., [4]),

$$R_{2/3}(f) \geq IC_{\mu, 2/3}(f),$$

where $R_{2/3}(f)$ is the communication complexity of the best randomized protocol (with error probability at most $2/3$) on the worst case input, μ is any distribution over the input space, and IC_μ is the information complexity of the best protocol over the input distribution μ . We now construct an input distribution μ as follows. Let D be chosen uniformly at random over $\{\text{Alice}, \text{Bob}\}$, E is chosen uniformly from $\{a, a + 1, a + 2, \dots, m - a\}$, where $m \gg a$ is a sufficiently large integer for which the coordinates of the vector jointly created by Alice and Bob will never exceed m . Denote by ξ the joint distribution of (D, E) . Based on the value of D and E , the input distribution γ of Alice and Bob is decided as follows: if D chooses Alice, then the input X of Alice is chosen uniformly at random from the multiset $\{E - a, E - b, E, E, E, E, E + a, E + b\}$, and Bob is given E ; if D chooses

Bob, γ just swaps the role of Alice and Bob. Clearly, γ is a product distribution conditioned on (D, E) . We have that $\zeta = (\gamma, \xi)$ is a mixture of product distributions. Let $\eta = \zeta^n$, and let μ be the distribution obtained by marginalizing (\mathbf{D}, \mathbf{E}) from η . Since every joint vector created above is from V_0 , μ is a *collapsing distribution* for f (see [4]). Then by Lemma 5.1 and Lemma 5.5 of [4], we have the following direct sum result:

$$I(\mathbf{X}, \mathbf{Y}; \Pi(\mathbf{X}, \mathbf{Y}) \mid \mathbf{D}, \mathbf{E}) \geq n \cdot \text{CIC}_{\zeta, 2/3}(\text{DIST}).$$

We now expand the terms of $\text{CIC}_{\zeta, 2/3}(\text{DIST})$. Let $[a, m - a] = \{a, a + 1, a + 2, \dots, m - a\}$, U_e be the random number from the uniform distribution on multiset $K \equiv \{e - a, e + a, e, e, e, e, e - b, e + b\}$, and Ψ be the transcript of a communication-optimal randomized protocol for the primitive function DIST . Let $\beta_{e,y}$ be the distribution of $\Psi(e, y)$ and $\beta_e = \frac{1}{8} \sum_{y \in K} \beta_{e,y}$ be the average distribution. The following inequalities follow from Lemma 2.45 and Lemma 2.52 in Bar-Yossef's PhD thesis [3], the Cauchy-Schwarz inequality, and the triangle inequality.

$$\begin{aligned} & \text{CIC}_{\zeta, 2/3}(\text{DIST}) \\ &= \frac{1}{2|m - 2a|} \sum_{e \in [a, m-a]} [I(U_e; \Psi(U_e, e)) \\ & \quad + I(U_e; \Psi(e, U_e))] \\ &= \sum_{\substack{e \in [a, m-a] \\ e_j \in K}} \frac{\text{Pr}[U_e = e_j]}{2(m - 2a)} [KL(\beta_{e, e_j} \parallel \beta_e) \\ & \quad + KL(\beta_{e_j, e} \parallel \beta_e)] \\ & \quad (\text{Bar-Yossef's Thesis Proposition 2.45}) \\ &\geq \frac{1}{8 \ln 2(m - 2a)} \sum_{\substack{e \in [a, m-a] \\ j \in \{-a, -b, 0, \dots, 0, a, b\}}} [h^2(\beta_{e+j, e}, \beta_e) \\ & \quad + h^2(\beta_{e, e+j}, \beta_e)] \\ & \quad (\text{Cauchy-Schwarz}) \\ &\geq \frac{1}{16 \ln 2(m - 2a)} \sum_{e \in [a, m-a]} [(h(\beta_{e,e}, \beta_e) + h(\beta_{e, e+a}, \beta_e))^2 \\ & \quad + (h(\beta_{e-a, e}, \beta_e) + h(\beta_{e, e}, \beta_e))^2 \\ & \quad + (h(\beta_{e, e}, \beta_e) + h(\beta_{e, e+b}, \beta_e))^2 \\ & \quad + (h(\beta_{e-b, e}, \beta_e) + h(\beta_{e, e}, \beta_e))^2] \\ & \quad (\text{Triangle inequality and Cauchy-Schwarz}) \\ &\geq \frac{C}{m} \sum_{e \in [2a, m-2a]} h^2(\beta_{e, e}, \beta_{e+a, e+a}) + h^2(\beta_{e, e}, \beta_{e+b, e+b}), \end{aligned} \tag{3}$$

where C is a constant. Now we are able to group the terms. By the Euclidean algorithm, there exist integers p and q such that $pa + qb = 1$. Let q be such an integer with smallest absolute value. We then have $a/b \leq |q| \leq a$ and $|p| \leq |q|$. Therefore, for any $i \in [4, m/4a - 4]$, we can select up to $|p| + |q|$ terms from the range $[4ai - 2a, 4ai + 2a]$ to group in the above equation. Without loss of generality, assume

$y > 0$. First, using the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \text{CIC}_{\zeta, 2/3}(\text{DIST}) &\geq \frac{C}{4ma} \sum_i \left(\sum_{e \in [4ai - 2a, 4ai + 2a]} \right. \\ & \quad h(\beta_{e, e}, \beta_{e+a, e+a}) + h(\beta_{e, e}, \beta_{e-a, e-a}) \\ & \quad \left. + h(\beta_{e, e}, \beta_{e+b, e+b}) + h(\beta_{e, e}, \beta_{e-b, e-b}) \right)^2. \end{aligned}$$

Next, group the terms using the triangle inequality such that whenever combining an a -term $h(\beta_{e, e}, \beta_{e'+a, e'+a})$, follow this by combining $\lfloor q/p \rfloor$ b -terms $h(\beta_{e, e}, \beta_{e''-b, e''-b})$. Therefore, the combined term e' is always guaranteed to be in $[4ai - 2a, 4ai + 2a]$. There might be a concern that a term is combined twice, namely, that there exist $|q_1| \leq |q_2| \leq |q|$, $|p_1| \leq |p_2| \leq |p|$ and $(q_1, p_1) \neq (q_2, p_2)$ such that $p_1 a + q_1 b = p_2 a + q_2 b$. But this is impossible since $\text{gcd}(a, b) = 1$ and $(p - (p_2 - p_1))a + (q - (q_2 - q_1))b = 1$ contradicts that q has the smallest absolute value. Therefore, we are left with,

$$\text{CIC}_{\zeta, 2/3}(\text{DIST}) \geq \frac{C'}{ma} \sum_{\substack{i \in [4, m/4a - 4] \\ e = 4ai}} h^2(\beta_{e, e}, \beta_{e+1, e+1}). \tag{4}$$

Now invoke the Pythagorean lemma of [4], stating that $h^2(\beta_{e, e}, \beta_{e+1, e+1}) \geq 1/2(h^2(\beta_{e, e}, \beta_{e+1, e}) + h^2(\beta_{e+1, e}, \beta_{e+1, e+1}))$, as well as the correctness of the protocol, stating that the $h^2(\beta_{e, e}, \beta_{e, e+1})$ terms can each be lower bounded by a positive constant. Thus,

$$\text{CIC}_{\zeta, 2/3}(\text{DIST}) \geq \frac{C''}{a^2}. \tag{5}$$

□

The above proof also makes use of the following lemma.

LEMMA 50. *Let a, b be two integers such that $\text{gcd}(a, b) = 1$ and $b < a$. Then there exist integers x, y such that $ax + by = 1$. Let y be such an integer with smallest absolute value. Then $b/a \leq |y| \leq a$ and $|x| \leq |y|$.*

PROOF. To see that $|y| \leq a$, w.l.o.g., we assume $y > a$. Then $y = qa + r$, where $r < a < y$ and $q < y$. Thus $ax + b(qa + r) = 1$, and $(qb + x)a + rb = 1$ contradicts that y has the smallest absolute value. It is clear that $|y| \geq a/b$. It remains to show that $|x| \leq |y|$. This holds since $x = (1 - by)/a$ and $b < a$. □

Now we look at a more general case.

THEOREM 51. *Let a, b, c be positive integers such that $c \neq a$ and $c \neq b$. Suppose there exist integers p, q such that $ap + bq = c$. Let q be such an integer with smallest absolute value. Then any randomized algorithm solving (a, b, c) - DIST with probability at least $2/3$ requires $\Omega(n/q^2)$ bits of memory.*

PROOF. The proof of this theorem is a modification of the proof of the previous proposition. With a similar distribution (replacing 1 with c), we have that Equation (3) holds. The remaining task is to show a partition scheme of the terms such that $\Theta(m/(|p| + |q|)^2)$ terms survive. Choose an arbitrary $e \in [4a^2, m - 4a^2]$. W.l.o.g. assume $p > 0$. We

can combine the following $|p| + |q|$ terms,

$$\begin{aligned}
& h^2(\beta_{e,e}, \beta_{e+a,e+a}) + h^2(\beta_{e+a,e+a}, \beta_{e+2a,e+2a}) \\
& + \dots + \\
& + h^2(\beta_{e+pa-a,e+pa-a}, \beta_{e+pa,e+pa}) \\
& + h^2(\beta_{e+pa,e+pa}, \beta_{e+pa-b,e+pa-b}) \\
& + h^2(\beta_{e+pa-b,e+pa-b}, \beta_{e+pa-2b,e+pa-2b}) \\
& + \dots + \\
& + h^2(\beta_{e+pa+qb+b,e+pa+qb+b}, \beta_{e+pa+qb,e+pa+qb}) \\
& \geq \frac{1}{|p| + |q|} \left[\sum_{i=1}^p h(\beta_{e+(i-1)a,e+(i-1)a}, \beta_{e+ia,e+ia}) \right. \\
& \left. + \sum_{j=1}^{|q|} h(\beta_{e+pa-(j-1)b,e+pa-(j-1)b}, \beta_{e+pa-jb,e+pa-jb}) \right]^2 \\
& \quad \text{(Cauchy-Schwarz)} \\
& \geq \frac{h^2(\beta_{e,e}, \beta_{e+1,e+1})}{|p| + |q|} \quad \text{(Triangle inequality)}.
\end{aligned}$$

Choose an e . This forbids the choice of another e from the above terms. In fact, e uniquely determines two arithmetic progressions: $e, e + a, e + 2a, \dots, e + pa$ and $e + pa - b, e + pa - 2b, \dots, e + pa + qb$. Now, in order to choose a new e' , the value e' cannot be a number from the above two progressions, and the progressions determined by e' cannot share a term with one of e . Therefore, because we chose e , $2(|p| + |q|)$ terms are not available for a new choice of e' . By choosing e from $[4a^2, m - 4a^2]$ one at a time, we can find $(m - 8a^2)/a$ different e for the purpose of combining terms. Denote the set of these e by S . In summary, we have,

$$\text{CIC}_{\zeta, 2/3}(\text{DIST}) \geq \frac{C}{m(|p| + |q|)} \sum_{e \in S} h^2(\beta_{e,e}, \beta_{e+1,e+1}). \quad (6)$$

Since $h^2(\beta_{e,e}, \beta_{e+1,e+1})$ is bounded below by a constant as indicated in the previous proposition, and since $|S| = m - 8a^2 = \Omega(m)$, we have $\text{CIC}_{\zeta, 2/3}(\text{DIST}) = \Omega(1/q^2)$. \square

The above lower bound is tight in the sense that there is a protocol that uses $O(n/q^2)$ memory bits and solves (a, b, c) -DIST.

PROPOSITION 52. *Let a, b, c, p, q be integers satisfying the same conditions as in the previous theorem. Then there is a randomized algorithm solving (a, b, c) -DIST with probability at least $2/3$ using $O(n/q^2) \cdot \text{poly}(\log n)$ bits of memory.*

PROOF. Before the beginning of the stream, we partition $[n]$ into t contiguous pieces each of size n/t , where $t = \tilde{O}(n/q^2)$, and the \tilde{O} notation hides logarithmic factors in n . For the i -th piece of the universe, let the corresponding substream restricted to elements in the i -th piece be denoted by S_i . For each S_i , we maintain a counter C_i for which

$$C_i = \sum_{l: h[l]=i} v_l \xi_l = \sum_{x \in \{a, b, 1\}} z_{i,x} x$$

where $\xi_l \sim \{-1, +1\}$ are uniform 4-wise independent random bits and $z_{i,x} = \sum_{l: h[l]=i, |v_l|=x} \xi_l$. Let $y_{i,x}$ denote the number of occurrence of $|x|$ in the i -th stream. We have $y_{i,x} \leq n/t$. By standard concentration arguments, we have that with arbitrarily large constant probability, $|z_{i,x}| = \tilde{O}(\sqrt{y_{i,x}}) = \tilde{O}(\sqrt{n/t})$. We can select an appropriate $t = \tilde{O}(n/q^2)$ for

which this implies $|z_{i,x}|$ is at most $|q|/4$. Now the claim is that by reading the value $C_i \bmod a$, this is enough to distinguish whether there is a frequency of absolute value “ c ” in the stream or not. This follows since the sets of values of $C_i \bmod a$ in the two cases are disjoint. To see why, note that if $z'_b b = z_b b + c \bmod a$, then we have $(z'_b - z_b) - ra = c$. However, since $|z'_b - z_b| < |q|$, by the minimality of $|q|$, this is a contradiction. \square

C.2 ShortLinearCombination Problem for Multiple of Frequencies

DEFINITION 53. *Let $u = (u_1, u_2, \dots, u_r)$ be a vector in \mathbb{Z}^r for an integer r . Let $d > 0$ be an integer not in the vector u . A stream S with frequency vector v is given to an algorithm, where v is promised to be from $V_0 = \{u_1, u_2, \dots, u_r, 0\}^n$ or $V_1 = \{\text{EMB}(v, i, e) \mid v \in V_0, i \in [n], e \in \{-d, d\}\}$. The (u, d) -DIST problem is to distinguish whether $v \in V_0$ or $v \in V_1$.*

THEOREM 54. *Let $u = (a_1, a_2, \dots, a_r)$ be a set of integers and $d > 0$ be a positive integer such that $d = q_1 a_1 + q_2 a_2 + \dots + q_r a_r$ where q_1, q_2, \dots, q_r are integers. These integers are chosen such that $q = \sum_i |q_i|$ is minimal. Then any randomized algorithm solving (u, d) -DIST with probability at least $2/3$ requires $\Omega(n/q^2)$ bits of memory.*

PROOF. The proof of the lower bound is a straightforward extension of the previous argument for three frequencies. Let $a = \max(a_1, a_2, \dots, a_r)$. The number of combined terms is $m/(|q_1| + |q_2| + \dots + |q_r|)$. Therefore, the lower bound is $\Omega(n/q^2)$. The upper bound is also a straightforward extension of the three frequency case. \square