

# A Framework for Estimating Stream Expression Cardinalities

Anirban Dasgupta\*

Kevin Lang<sup>†</sup>Lee Rhodes<sup>‡</sup>Justin Thaler<sup>§</sup>

## Abstract

Given  $m$  distributed data streams  $A_1, \dots, A_m$ , we consider the problem of estimating the number of unique identifiers in streams defined by set expressions over  $A_1, \dots, A_m$ . We identify a broad class of algorithms for solving this problem, and show that the estimators output by any algorithm in this class are perfectly unbiased and satisfy strong variance bounds. Our analysis unifies and generalizes a variety of earlier results in the literature. To demonstrate its generality, we describe several novel sampling algorithms in our class, and show that they achieve a novel tradeoff between accuracy, space usage, update speed, and applicability.

## 1 Introduction

Consider an internet company that monitors the traffic flowing over its network by placing a sensor at each ingress and egress point. Because the volume of traffic is large, each sensor stores only a small *sample* of the observed traffic, using some simple sampling procedure. At some later point, the company decides that it wishes to estimate the number of unique users who satisfy a certain property  $P$  and have communicated over its network. We refer to this as the  $\text{DISTINCTONSUBPOPULATION}_P$  problem, or  $\text{DISTINCT}_P$  for short. How can the company combine the samples computed by each sensor, in order to accurately estimate the answer to this query?

In the case that  $P$  is the trivial property that is satisfied by all users, the answer to the query is simply the number of  $\text{DISTINCTELEMENTS}$  in the traffic stream, or  $\text{DISTINCT}$  for short. The problem of designing streaming algorithms and sampling procedures for estimating  $\text{DISTINCTELEMENTS}$  has been the subject of intense study. In general, however,  $P$  may be significantly more complicated than the trivial property, and may not be known until query time. For example, the company may want to estimate the number of (unique) men in a certain age range, from a specified country, who accessed a certain set of websites during a designated time period, while excluding IP addresses belonging to a designated blacklist. This more general setting, where  $P$  is a nontrivial ad hoc property, has received somewhat less attention than the basic  $\text{DISTINCT}$  problem.

In this paper, our goal is to identify a simple method for combining the samples from each sensor, so that the following holds. As long as each sensor is using a sampling procedure that satisfies a certain mild technical condition, then for any property  $P$ , the combining procedure outputs an estimate for the  $\text{DISTINCT}_P$  problem that is unbiased. Moreover, its variance should be bounded by that of the individual sensors' sampling procedures.<sup>1</sup>

For reasons that will become clear later, we refer to our proposed combining procedure as the *Theta-Sketch Framework*, and we refer to the mild technical condition that each sampling procedure must satisfy to guarantee unbiasedness as *1-Goodness*. If the sampling procedures satisfy an additional property that we refer to as *monotonicity*, then the variance of the estimate output by the combining procedure is guaranteed to satisfy the desired variance bound. The Theta-Sketch Framework, and our analysis of it, unifies and generalizes a variety of results in the literature (see Section 2.5 for details).

---

\*Indian Institute of Technology, Gandhinagar

<sup>†</sup>Yahoo Labs

<sup>‡</sup>Yahoo! Inc

<sup>§</sup>Yahoo Labs

<sup>1</sup>More precisely, we are interested in showing that the variance of the returned estimate is at most that of the (hypothetical) estimator obtained by running each individual sensor's sampling algorithm on the concatenated stream  $A_1 \circ \dots \circ A_m$ . We refer to the latter estimator as "hypothetical" because it is typically infeasible to materialize the concatenated stream in distributed environments.

**The Importance of Generality.** As we will see, there is a huge array of sampling procedures that the sensors could use. Each procedure comes with a unique tradeoff between accuracy, space requirements, update speed, and simplicity. Moreover, some of these procedures come with additional desirable properties, while others do not. We would like to support as many sampling procedures as possible, because the best one to use in any given setting will depend on the relative importance of each resource in that setting.

**Handling Set Expressions.** The scenario described above can be modeled as follows. Each sensor observes a stream of identifiers  $A_j$  from a data universe of size  $n$ , and the goal is to estimate the number of distinct identifiers that satisfy property  $P$  in the combined stream  $U = \cup_j A_j$ . In full generality, we may wish to handle more complicated set expressions applied to the constituent streams, other than set-union. For example, we may have  $m$  streams of identifiers  $A_1, \dots, A_m$ , and wish to estimate the number of distinct identifiers satisfying property  $P$  that appear in *all streams*. The Theta-Sketch Framework can be naturally extended to provide estimates for such queries. Our analysis applies to any sequence of set operations on the  $A_j$ 's, but we restrict our attention to set-union and set-intersection throughout the paper for simplicity.

## 2 Preliminaries, Background, and Contributions

### 2.1 Notation and Assumptions

**Streams and Set Operations.** Throughout,  $A$  denotes a stream of identifiers from a data universe  $[n] := \{1, \dots, n\}$ . We view any *property*  $P$  on identifiers as a subset of  $[n]$ , and let  $n_{P,A} := \text{DISTINCT}_P(A)$  denote the number of distinct identifiers that appear in  $A$  and satisfy  $P$ . For brevity, we let  $n_A$  denote  $\text{DISTINCT}(A)$ . When working in a multi-stream setting,  $A_1, \dots, A_m$  denote  $m$  streams of identifiers from  $[n]$ ,  $U := \cup_{j=1}^m A_j$  will denote the concatenation of the  $m$  input streams, while  $I := \cap_{j=1}^m A_j$  denotes the set of identifiers that appear at least once in all  $m$  streams. Because we are interested only in *distinct* counts, it does not matter for definitional purposes whether we view  $U$  and  $I$  as sets, or as multisets. For any property  $P: [n] \rightarrow \{0, 1\}$ ,  $n_{P,U} := \text{DISTINCT}_P(U)$  and  $n_{P,I} := \text{DISTINCT}_P(I)$ , while  $n_U := \text{DISTINCT}(U)$  and  $n_I := \text{DISTINCT}(I)$ .

**Hash Functions.** For simplicity and clarity, and following prior work (e.g. [5, 9]), we assume throughout that the sketching and sampling algorithms make use of a perfectly random hash function  $h$  mapping the data universe  $[n]$  to the open interval  $(0, 1)$ . That is, for each  $x \in [n]$ ,  $h(x)$  is a uniform random number in  $(0, 1)$ . Given a subset of hash values  $S$  computed from a stream  $A$ , and a property  $P \subseteq [n]$ ,  $P(S)$  denotes the subset of hash values in  $S$  whose corresponding identifiers in  $[n]$  satisfy  $P$ . Finally, given a stream  $A$ , the notation  $X^{n,A}$  refers to the set of hash values obtained by mapping a hash function  $h$  over the  $n_A$  distinct identifiers in  $A$ .

### 2.2 Prior Art: Sketching Procedures for DISTINCT Queries

There is a sizeable literature on streaming algorithms for estimating the number of distinct elements in a single data stream. Some, but not all, of these algorithms can be modified to solve the  $\text{DISTINCT}_P$  problem for general properties  $P$ . Depending on which functionality is required, systems based on HyperLogLog Sketches, K'th Minimum Value (KMV) Sketches, and Adaptive Sampling represent the state of the art for practical systems [21].<sup>2</sup> For clarity of exposition, we defer a thorough overview of these algorithms to Section 6. Here, we briefly review the main concepts and relevant properties of each.

**HLL: HyperLogLog Sketches.** HLL is a sketching algorithm for the vanilla  $\text{DISTINCT}$  problem. Its accuracy per bit is superior to the KMV and Adaptive Sampling algorithms described below. However, unlike KMV and Adaptive Sampling, it is not known how to extend the HLL sketch to estimate  $n_{P,A}$  for general properties  $P$  (unless, of course,  $P$  is known prior to stream processing).

**KMV: K'th Minimum Value Sketches.** The KMV sketching procedure for estimating  $\text{DISTINCT}(A)$  works as follows. While processing an input stream  $A$ , KMV keeps track of the set  $S$  of the  $k$  smallest unique hashed values of stream

<sup>2</sup>Algorithms with better asymptotic bit-complexity are known [23], but they do not match the practical performance of the algorithms discussed here. See Section 6.3.

elements. The update time of a heap-based implementation of KMV is  $O(\log k)$ . The KMV estimator for  $\text{DISTINCT}(A)$  is:  $\text{KMV}_A = k/m_{k+1}$ , where  $m_{k+1}$  denotes the  $k+1^{\text{st}}$  smallest unique hash value.<sup>3</sup> It has been proved by [5], [19], and others, that  $E(\text{KMV}_A) = n_A$ , and  $\sigma^2(\text{KMV}_A) = \frac{n_A^2 - k n_A}{k-1} < \frac{n_A^2}{k-1}$ . Duffield et al. [11] proposed to change the heap-based implementation of the KMV sketching algorithm to an implementation based on quickselect [22]. This reduces the sketch update cost from  $O(\log k)$  to amortized  $O(1)$ . However, this  $O(1)$  hides a larger constant than competing methods. At the cost of storing the sampled identifiers, and not just their hash values, the KMV sketching procedure can be extended to estimate  $n_{P,A}$  for any property  $P \subseteq [n]$  (Section 6 has details).

**Adaptive Sampling.** Adaptive Sampling maintains a sampling level  $i \geq 0$ , and the set  $S$  of all hash values less than  $2^{-i}$ ; whenever  $|S|$  exceeds a pre-specified size limit,  $i$  is incremented and  $S$  is scanned discarding any hash value that is now too big. Because a simple scan is cheaper than running quickselect, an implementation of this scheme is typically faster than KMV. The estimator of  $n_A$  is  $\text{Adapt}_A = |S|/2^{-i}$ . It has been proved by [13] that this estimator is unbiased, and that  $\sigma^2(\text{Adapt}_A) \approx 1.44(n_A^2/(k-1))$ , where the approximation sign hides oscillations caused by the periodic culling of  $S$ . Like KMV, Adaptive Sampling can be extended to estimate  $n_{P,A}$  for any property  $P$ . Although the stream processing speed of Adaptive Sampling is excellent, the fact that its accuracy oscillates as  $n_A$  increases is a shortcoming.

**HLL for set operations on streams.** HLL can be directly adapted to handle set-union (see Section 6 for details). For set-intersection, the relevant adaptation uses the inclusion/exclusion principle. However, the variance of this estimate is approximately a factor of  $n_U/n_I$  worse than the variance achieved by the multiKMV algorithm described below. When  $n_I \ll n_U$ , this penalty factor overwhelms HLL’s fundamentally good accuracy per bit.

**KMV for set operations on streams.** Given streams  $A_1, \dots, A_m$ , let  $S_j$  denote the KMV sketch computed from stream  $A_j$ . A trivial way to use these sketches to estimate the number of distinct items  $n_U$  in the union stream  $U$  is to let  $M'_U$  denote the  $(k+1)^{\text{st}}$  smallest value in the union of the sketches, and let  $S'_U = \{x \in \cup_j S_j : x < M'_U\}$ . Then  $S'_U$  is identical to the sketch that would have been obtained by running KMV directly on the concatenated stream  $A_1 \circ \dots \circ A_m$ , and hence  $\text{KMV}_{P,U} := k/M'_U$  is an unbiased estimator for  $n_U$ , by the same analysis as in the single-stream setting. We refer to this procedure as the “non-growing union rule.”

Intuitively, the non-growing union rule does not use all of the information available to it. The sets  $S_j$  contain up to  $k \cdot M$  distinct samples in total, but  $S'_U$  ignores all but the  $k$  smallest samples. With this in mind, Cohen and Kaplan [9] proposed the following adaptation of KMV to handle unions of multiple streams. We denote their algorithm by multiKMV, and also refer to it as the “growing union rule”.

For each KMV sketch  $S_j$  computed from stream  $A_j$ , let  $M_j$  denote that sketch’s value of  $m_{k+1}$ . Define  $M_U = \min_{j=1}^m M_j$ , and  $S_U = \{x \in \cup_j S_j : x < M_U\}$ . Then  $n_U$  is estimated by  $\text{multiKMV}_U := |S_U|/M_U$ , and  $n_{P,U}$  is estimated by  $\text{multiKMV}_{P,U} := |P(S_U)|/M_U$ .

At first glance, it may seem obvious that the growing union rule yields an estimator that is “at least as good” as the non-growing union, since the growing union rule makes use of at least as many samples as the non-growing rule. However, it is by no means trivial to prove that  $\text{multiKMV}_{P,U}$  is unbiased, nor that its variance is dominated by that of the non-growing union rule. Nonetheless, [9] managed to prove this: they showed that  $\text{multiKMV}_{P,U}$  is unbiased and has variance that is dominated by the variance of  $\text{KMV}_{P,U}$ :

$$\sigma^2(\text{multiKMV}_{P,U}) \leq \sigma^2(\text{KMV}_{P,U}). \quad (1)$$

As observed in [9], multiKMV can be adapted in a similar manner to handle set-intersections (see Section 3.8 for details).

**Adaptive Sampling for set operations on streams.** Adaptive Sampling can handle set unions and intersections with a similar “growing union rule” in which “ $M_U$ ” :=  $\min_{j=1}^m (2^{-i})_j$ . Here,  $(2^{-i})_j$  denotes the threshold for discarding hash values that was computed by the  $j$ th Adaptive Sampling sketch. We refer to this algorithm as multiAdapt. [18] proved epsilon-delta bounds on the error of  $\text{multiAdapt}_{P,U}$ , but did not derive expressions for mean or variance. However, multiAdapt and multiKMV are both special cases of our Theta-Sketch Framework, and in Section 3 we will prove (apparently for the first time) that  $\text{multiAdapt}_{P,U}$  is unbiased, and satisfies strong variance bounds. These results

<sup>3</sup>Some works use the estimate  $k/m_k$ , e.g. [4]. We use  $k/m_{k+1}$  because it is unbiased, and for consistency with the work of Cohen and Kaplan [9] described below.

---

**Algorithm 1** Theta Sketch Framework for estimating  $n_{P,U}$ . The framework is parameterized by choice of TCF's  $T^{(j)}(k, A_j, h)$ , one for each input stream.

---

- 1: **Definition:** Function  $\text{samp}_j[T^{(j)}](k, A_j, h)$
  - 2:  $\theta_j \leftarrow T^{(j)}(k, A_j, h)$
  - 3:  $S_j \leftarrow \{(x \in h(A_j)) < \theta_j\}$ .
  - 4: **return**  $(\theta_j, S_j)$ .
  - 5: **Definition:** Function  $\text{ThetaUnion}(\text{Theta Sketches } \{(\theta_j, S_j)\})$
  - 6:  $\theta_U \leftarrow \min\{\theta_j\}$ .
  - 7:  $S_U \leftarrow \{(x \in (\cup S_j)) < \theta_U\}$ .
  - 8: **return**  $(\theta_U, S_U)$ .
  - 9: **Definition:** Function  $\text{EstimateOnSubPopulation}(\text{Theta Sketch } (\theta, S)$  produced from stream  $A$ , Property  $P$  mapping identifiers to  $\{0, 1\}$ )
  - 10: **return**  $\hat{n}_{A,P} := \frac{|P(S)|}{\theta}$ .
- 

have the following two advantages over the epsilon-delta bounds of [18]. First, proving unbiasedness is crucial for obtaining estimators for distinct counts over subpopulations: these estimators are analyzed as a sum of a huge number of per-item estimates (see Theorem 3.10 for details), and biases add up. Second, variance bounds enable derivation of confidence intervals that an epsilon-delta guarantee cannot provide, unless the guarantee holds for many values of delta simultaneously.

### 2.3 Overview of the Theta-Sketch Framework

In this overview, we describe the Theta-Sketch Framework in the multi-stream setting where the goal is to output  $n_{P,U}$ , where  $U = \cup_{j=1}^m A_j$  (we define the framework formally in Section 2.4). That is, the goal is to identify a very large class of sampling algorithms that can run on each constituent stream  $A_j$ , as well as a “universal” method for combining the samples from each  $A_j$  to obtain a good estimator for  $n_{P,U}$ . We clarify that the Theta-Sketch Framework, and our analysis of it, yields unbiased estimators that are interesting even in the single-stream case, where  $m = 1$ .

We begin by noting the striking similarities between the multiKMV and multiAdapt algorithms outlined in Section 2.2. In both cases, a sketch can be viewed as pair  $(\theta, S)$  where  $\theta$  is a certain threshold that depends on the stream, and  $S$  is a set of hash values which are all strictly less than  $\theta$ . In this view, both schemes use the same estimator  $|S|/\theta$ , and also the same growing union rule for combining samples from multiple streams. The only difference lies in their respective rules for mapping streams to thresholds  $\theta$ . The Theta-Sketch Framework formalizes this pattern of similarities and differences.

**The assumed form of the single-stream sampling algorithms.** The Theta-Sketch Framework demands that each constituent stream  $A_j$  be processed by a sampling algorithm  $\text{samp}_j$  of the following form. While processing  $A_j$ ,  $\text{samp}_j$  evaluates a “threshold choosing function” (TCF)  $T^{(j)}(A_j)$ . The final state of  $\text{samp}_j$  must be of the form  $(\theta_j := T^{(j)}(A_j), S)$ , where  $S$  is the set of all hash values strictly less than  $\theta_j$  that were observed while processing  $A_j$ . If we want to estimate  $n_{P,U}$  for non-trivial properties  $P$ , then  $\text{samp}_j$  must also store the corresponding identifier that hashed to each value in  $S$ . Note that the framework itself does not specify the threshold-choosing functions  $T^{(j)}$ . Rather, any specification of the TCFs  $T^{(j)}$  defines a particular instantiation of the framework.

**Remark.** It might appear from Algorithm 1 that for any TCF  $T^{(j)}$ , the function  $\text{samp}_j[T^{(j)}]$  makes two passes over the input stream: one to compute  $\theta_j$ , and another to compute  $S_j$ . However, in all of the instantiations we consider, both operations can be performed in a single pass.

**The universal combining rule.** Given the states  $(\theta_j := T^{(j)}(A_j), S_j)$  of each of the  $m$  sampling algorithms when run on the streams  $A_1, \dots, A_m$ , define  $\theta_U := \min_{j=1}^m \theta_j$ , and  $S_U := \{x \in \cup_j S_j : x < \theta_U\}$  (see the function  $\text{ThetaUnion}$  in Algorithm 1). Then  $n_U$  is estimated by  $\hat{n}_U := |S_U|/\theta_U$ , and  $n_{P,U}$  as  $\hat{n}_{P,U} := |P(S_U)|/\theta_U$  (see the function  $\text{EstimateOnSubPopulation}$  in Algorithm 1).

**The analysis.** Our analysis shows that, so long as each threshold-choosing function  $T^{(j)}$  satisfies a mild technical condition that we call *I-Goodness*, then  $\hat{n}_{P,U}$  is unbiased. We also show that if each  $T^{(j)}$  satisfies a certain additional condition that we call *monotonicity*, then  $\hat{n}_{P,U}$  satisfies strong variance bounds (analogous to the bound of Equation

(1) for KMV). Our analysis is arguably surprising, because 1-Goodness does not imply certain properties that have traditionally been considered important, such as permutation invariance, or  $S$  being a uniform random sample of the hashed unique items of the input stream.

**Applicability.** To demonstrate the generality of our analysis, we identify several valid instantiations of the Theta-Sketch Framework. First, we show that the TCF’s used in KMV and Adaptive Sampling both satisfy 1-Goodness and monotonicity, implying that multiKMV and multiAdapt are both unbiased and satisfy the aforementioned variance bounds. For multiKMV, this is a reproof of Cohen and Kaplan’s results [9], but for multiAdapt the results are new. Second, we identify a variant of KMV that we call pKMV, which is useful in multi-stream settings where the lengths of constituent streams are highly skewed. We show that pKMV satisfies both 1-Goodness and monotonicity. Third, we introduce a new sampling procedure that we call the *Alpha Algorithm*. Unlike earlier algorithms, the Alpha Algorithm’s final state actually depends on the stream order, yet we show that it satisfies 1-Goodness, and hence is unbiased in both the single- and multi-stream settings. We also establish variance bounds on the Alpha Algorithm in the single-stream setting. We show experimentally that the Alpha Algorithm, in both the single- and multi-stream settings, achieves a novel tradeoff between accuracy, space usage, update speed, and applicability.

Unlike KMV and Adaptive Sampling, the Alpha Algorithm does not satisfy monotonicity in general. In fact, we have identified contrived examples in the multi-stream setting on which the aforementioned variance bounds are (weakly) violated. The Alpha Algorithm does, however, satisfy monotonicity under the promise that the  $A_1, \dots, A_m$  are pairwise disjoint, implying variance bounds in this case. Our experiments suggest that, in practice, the normalized variance in the multi-stream setting is not much larger than in the pairwise disjoint case.

**Deployment of Algorithms.** Within Yahoo, the pKMV and Alpha algorithms are used widely. In particular, stream cardinalities in Yahoo empirically satisfy a power law, with some very large streams and many short ones, and pKMV is an attractive option for such settings. We have released an optimized open-source implementation of our algorithms at <http://datasketches.github.io/>.

## 2.4 Formal Definition of Theta-Sketch Framework

The Theta-Sketch Framework is defined as follows. This definition is specific to the multi-stream setting where the goal is to output  $n_{P,U}$ , where  $U = \cup_{j=1}^m A_j$  is the union of constituent streams  $A_1, \dots, A_m$ .

**Definition 2.1.** *The Theta-Sketch Framework consists of the following components:*

- *The data type  $(\theta, S)$ , where  $0 < \theta \leq 1$  is a threshold, and  $S$  is the set of all unique hashed stream items  $0 \leq x < 1$  that are less than  $\theta$ . We will generically use the term “theta-sketch” to refer to an instance of this data type.*
- *The universal “combining function”  $\text{ThetaUnion}()$ , defined in Algorithm 1, that takes as input a collection of theta-sketches (purportedly obtained by running  $\text{samp}[T]()$  on constituent streams  $A_1, \dots, A_m$ ), and returns a single theta-sketch (purportedly of the union stream  $U = \cup_{i=1}^m A_i$ ).*
- *The function  $\text{EstimateOnSubPopulation}()$ , defined in Algorithm 1, that takes as input a theta-sketch  $(\theta, S)$  (purportedly obtained from some stream  $A$ ) and a property  $P \subseteq [n]$  and returns an estimate of  $\hat{n}_{P,A}$ .*

Any instantiation of the Theta-Sketch Framework must specify a “threshold choosing function” (TCF), denoted  $T(k, A, h)$ , that maps a target sketch size, a stream, and a hash function  $h$  to a threshold  $\theta$ . Any TCF  $T$  implies a “base” sampling procedure  $\text{samp}[T]()$  that maps a target size, a stream  $A$ , and a hash function to a theta-sketch using the pseudocode shown in Algorithm 1. One can obtain an estimate  $\hat{n}_{P,A}$  for  $n_{P,A}$  by feeding the resulting theta-sketch into  $\text{EstimateOnSubPopulation}()$ .

Given constituent streams  $A_1, \dots, A_m$ , the instantiation obtains an estimate  $\hat{n}_{P,U}$  of  $n_{P,U}$  by running  $\text{samp}[T]()$  on each constituent stream  $A_j$ , feeding the resulting theta-sketches to  $\text{ThetaUnion}()$  to obtain a “combined” theta-sketch for  $U = \cup_{i=1}^m A_i$ , and then running  $\text{EstimateOnSubPopulation}()$  on this combined sketch.

**Remark.** Definition 2.1 assumes for simplicity that the same TCF  $T$  is used in the base sampling algorithms run on each of the constituent streams. However, all of our results that depend only on 1-Goodness (*e.g.* unbiasedness of estimates and non-correlation of “per-item estimates”) hold even if different 1-Good TCF’s are used on each stream, and even if different values of  $k$  are employed.

## 2.5 Summary of Contributions

In summary, our contributions are: (1) Formulating the Theta-Sketch Framework. (2) Identifying a mild technical condition (1-Goodness) on TCF’s ensuring that the framework’s estimators are unbiased. If each TCF also satisfies a monotonicity condition, the framework’s estimators come with strong variance bounds analogous to Equation (1). (3) Proving multiKMV, multiAdapt, and pKMV all satisfy 1-Goodness and monotonicity, implying unbiasedness and variance bounds for each. (4) Introducing the Alpha Algorithm, proving that it is unbiased, and establishing quantitative bounds on its variance in the single-stream setting. (5) Experimental results showing that the Alpha Algorithm instantiation achieves a novel tradeoff between accuracy, space usage, update speed, and applicability.

## 3 Analysis of the Theta-Sketch Framework

**Section Outline.** Section 3.1 shows that KMV and Adaptive Sampling are both instantiations of the Theta-Sketch Framework. Section 3.2 defines 1-Goodness. Sections 3.3 and 3.4 prove that the TCF’s that instantiate behavior identical to KMV and Adapt both satisfy 1-Goodness. Section 3.5 proves that if a framework instantiation’s TCF satisfies 1-Goodness, then so does the TCF that is implicitly applied to the union stream via the composition of the instantiation’s base algorithm and the function ThetaUnion(). Section 3.6 proves that the estimator  $\hat{n}_{P,A}$  for  $n_{P,A}$  returned by EstimateOnSubPopulation() is unbiased when applied to any theta-sketch produced by a TCF satisfying 1-Goodness. Section 3.7 defines monotonicity and shows that 1-Goodness and monotonicity together imply variance bounds on  $\hat{n}_{P,U}$ . Section 3.8 explains how to tweak the Theta-Sketch Framework to handle set intersections and other set operations on streams. Finally, Section 3.9 describes the pKMV variant of KMV.

### 3.1 Example Instantiations

Define  $m_{k+1}$  to be the  $k+1^{\text{st}}$  smallest unique hash value in  $h(A)$  (the hashed version of the input stream). The following is an easy observation.

**Observation 3.1.** *When the Theta-Sketch Framework is instantiated with the TCF  $T(k, A, h) = m_{k+1}$ , the resulting instantiation is equivalent to the multiKMV algorithm outlined in Section 2.2.*

Let  $\beta$  be any real value in  $(0, 1)$ . For any  $z$ , define  $\beta^{i(z)}$  to be the largest value of  $\beta^i$  (with  $i$  a non-negative integer) that is less than  $z$ .

**Observation 3.2.** *When the Theta-Sketch Framework is instantiated with the TCF  $T(k, A, h) = \beta^{i(m_{k+1})}$  the resulting instantiation is equivalent to multiAdapt, which combines Adaptive Sampling with a growing union rule (cf. Section 2.2).<sup>4</sup>*

### 3.2 Definition of 1-Goodness

The following circularity is a main source of technical difficulty in analyzing theta sketches: for any given identifier  $\ell$  in a stream  $A$ , whether its hashed value  $x_\ell = h(\ell)$  will end up in a sketch’s sample set  $S$  depends on a comparison of  $x_\ell$  versus a threshold  $T(X^{n_A})$  that depends on  $x_\ell$  itself. Adapting a technique from [9], we partially break this circularity by analyzing the following infinite family of projections of a given threshold choosing function  $T(X^{n_A})$ .

<sup>4</sup> Section 2.2 assumed that the parameter  $\beta$  was set to the most common value:  $1/2$ .

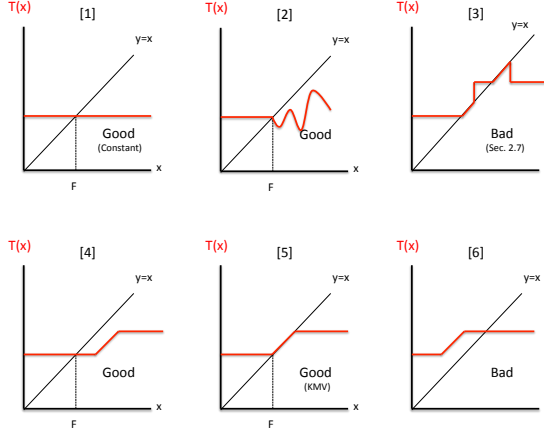


Figure 1: Six examples of hypothetical projections of TCF's. Four of them satisfy 1-Goodness; the other two do not.

**Definition 3.3** (Definition of Fix-All-But-One Projection). *Let  $T$  be a threshold choosing function. Let  $\ell$  be one of the  $n_A$  unique identifiers in a stream  $A$ . Let  $X_{-\ell}^{n_A}$  be a fixed assignment of hash values to all unique identifiers in  $A$  except for  $\ell$ . Then the fix-all-but-one projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell) : (0, 1) \rightarrow (0, 1)$  of  $T$  is the function that maps values of  $x_\ell$  to theta-sketch thresholds via the definition  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = T(X^{n_A})$ , where  $X^{n_A}$  is the obvious combination of  $X_{-\ell}^{n_A}$  and  $x_\ell$ .*

[9] analyzed similar projections under the assumption that the base algorithm is specifically (a weighted version of) KMV; we will instead impose the weaker condition that every fix-all-but-one projection satisfies 1-Goodness, defined below.<sup>5</sup>

**Definition 3.4** (Definition of 1-Goodness for Univariate Functions). *A function  $f(x) : (0, 1) \rightarrow (0, 1)$  satisfies 1-Goodness iff there exists a fixed threshold  $F$  such that:*

$$\text{If } x < F, \text{ then } f(x) = F. \quad (2)$$

$$\text{If } x \geq F, \text{ then } f(x) \leq x. \quad (3)$$

Figure 1 contains six examples of hypothetical projections of TCF's. Four of them satisfy 1-Goodness; the other two do not.

**Condition 3.5** (Definition of 1-Goodness for TCF's). *A TCF  $T(X^{n_A})$  satisfies 1-Goodness iff for every stream  $A$  containing  $n_A$  unique identifiers, every label  $\ell \in A$ , and every fixed assignment  $X_{-\ell}^{n_A}$  of hash values to the identifiers in  $A \setminus \ell$ , the fix-all-but-one projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  satisfies Definition 3.4.*

### 3.3 TCF of multiKMV Satisfies 1-Goodness

The following theorem shows that the TCF used in KMV satisfies 1-Goodness.

**Theorem 3.6.** *If  $T(X^{n_A}) = m_{k+1}$ , then every fix-all-but-one projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  of  $T$  satisfies 1-Goodness.*

*Proof.* Let  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  be any specific fix-all-but-one-projection of  $T(X^{n_A}) = m_{k+1}$ . We will exhibit the fixed value  $F_\ell[X_{-\ell}^{n_A}]$  that causes (2) and (3) to be true for this projection. Let  $a$  and  $b$  respectively be the  $k$ 'th and  $(k+1)$ 'st smallest hash values in  $X_{-\ell}^{n_A}$ . Then Subconditions (2) and (3) hold for  $F_\ell[X_{-\ell}^{n_A}] = a$ . There are three cases:

<sup>5</sup>We chose the name 1-Goodness due to the reference to Fix-All-But-One Projections.

**Case** ( $x_\ell < a < b$ ): In this case,  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = T(X^{n_A}) = m_{k+1} = a$ . Since  $x_\ell < (F_\ell[X_{-\ell}^{n_A}] = a)$ , (2) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = a) = F_\ell[X_{-\ell}^{n_A}]$ , and (3) holds vacuously.

**Case** ( $a < x_\ell < b$ ): In this case,  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = T(X^{n_A}) = m_{k+1} = x_\ell$ . Since  $x_\ell \geq (F_\ell[X_{-\ell}^{n_A}] = a)$ , (3) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = x_\ell) \leq x_\ell$ , and (2) holds vacuously.

**Case** ( $a < b < x_\ell$ ): In this case,  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = T(X^{n_A}) = m_{k+1} = b$ . Since  $x_\ell \geq (F_\ell[X_{-\ell}^{n_A}] = a)$ , (3) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = b) < x_\ell$ , and (2) holds vacuously. □

### 3.4 TCF of multiAdapt Satisfies 1-Goodness

The following theorem shows that the TCF used in Adaptive Sampling satisfies 1-Goodness.

**Theorem 3.7.** *If  $T(X^{n_A}) = \beta^{i(m_{k+1})}$ , then every fix-all-but-one projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  of  $T$  satisfies 1-Goodness.*

*Proof.* Let  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  be any specific fix-all-but-one-projection of  $T(X^{n_A}) = \beta^{i(m_{k+1})}$ . We will exhibit the fixed value  $F_\ell[X_{-\ell}^{n_A}]$  that causes (2) and (3) to be true for this projection. Let  $a$  and  $b$  respectively be the  $k$ 'th and  $(k+1)$ 'st smallest hash values in  $X_{-\ell}^{n_A}$ . Then Subconditions (2) and (3) hold for  $F_\ell[X_{-\ell}^{n_A}] = \beta^{i(a)}$ . There are four cases:

**Case** ( $x_\ell < \beta^{i(a)} < a < b$ ):  $m_{k+1} = a$ , so  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(a)}$ . Since  $x_\ell < (F_\ell[X_{-\ell}^{n_A}] = \beta^{i(a)})$ , (2) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(a)}) = F_\ell[X_{-\ell}^{n_A}]$ , and (3) holds vacuously.

**Case** ( $\beta^{i(a)} < x_\ell < a < b$ ):  $m_{k+1} = a$ , so  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(a)}$ . Since  $x_\ell \geq (F_\ell[X_{-\ell}^{n_A}] = \beta^{i(a)})$ , (3) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(a)}) < x_\ell$ , and (2) holds vacuously.

**Case** ( $\beta^{i(a)} < a < x_\ell < b$ ):  $m_{k+1} = x_\ell$ , so  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(x_\ell)}$ . Since  $x_\ell \geq (F_\ell[X_{-\ell}^{n_A}] = \beta^{i(a)})$ , (3) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(x_\ell)}) < x_\ell$ , and (2) holds vacuously.

**Case** ( $\beta^{i(a)} < a < b < x_\ell$ ):  $m_{k+1} = b$ , so  $T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(b)}$ . Since  $x_\ell \geq (F_\ell[X_{-\ell}^{n_A}] = \beta^{i(a)})$ , (3) holds because  $(T_\ell[X_{-\ell}^{n_A}](x_\ell) = \beta^{i(b)}) < b < x_\ell$ , and (2) holds vacuously. □

### 3.5 1-Goodness Is Preserved by the Function ThetaUnion()

Next, we show that if a framework instantiation's TCF  $T$  satisfies 1-Goodness, then so does the TCF  $T^U$  that is implicitly being used by the theta-sketch construction algorithm defined by the composition of the instantiation's base sampling algorithms and the function ThetaUnion(). We begin by formally extending the definition of a fix-all-but-one projection to cover the degenerate case where the label  $\ell$  isn't actually a member of the given stream  $A$ .

**Definition 3.8.** *Let  $A$  be a stream containing  $n_A$  identifiers. Let  $\ell$  be a label that is not a member of  $A$ . Let the notation  $X_{-\ell}^{n_A}$  refer to an assignment of hash value to all identifiers in  $A$ . For any hash value  $x_\ell$  of the non-member label  $\ell$ , define the value of the "fix-all-but-one" projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  to be the constant  $T(X_{-\ell}^{n_A})$ .*

**Theorem 3.9.** *If the threshold choosing functions  $T^{(j)}(X^{n_{A_j}})$  of the base algorithms used to create sketches of  $m$  streams  $A_j$  all satisfy Condition 3.5, then so does the TCF:*

$$T^U(X^{n_U}) = \min_j \{T^{(j)}(X^{n_{A_j}})\} \quad (4)$$

*that is implicitly applied to the union stream via the composition of those base algorithms and the procedure ThetaUnion().*

*Proof.* Let  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$  be any specific fix-all-but-one projection of the threshold choosing function  $T^U(X^{n_U})$  defined by Equation (4). We will exhibit the fixed value  $F^U[X_{-\ell}^{n_U}]$  that causes (2) and (3) to be true for  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ .



The projection  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$  is specified by a label  $\ell \in (A_U = \cup_j A_j)$ , and a set  $X_{-\ell}^{n_U}$  of fixed hash values for the identifiers in  $A_U \setminus \ell$ . For each  $j$ , those fixed hash values  $X_{-\ell}^{n_U}$  induce a set  $X_{-\ell}^{n_{A_j}}$  of fixed hash values for the identifiers in  $A_j \setminus \ell$ . The combination of  $\ell$  and  $X_{-\ell}^{n_{A_j}}$  then specifies a projection  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$  of  $T^{(j)}(X^j)$ . Now, if  $\ell \in A_j$ , this is a fix-all-but-one projection according to the original Definition 3.3, and according to the current theorem's pre-condition, this projection must satisfy 1-Goodness for univariate functions. On the other hand, if  $\ell \notin A_j$ , this is a fix-all-but-one projection according to the extended Definition 3.8, and is therefore a constant function, and therefore satisfies 1-Goodness. Because the projection  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$  satisfies 1-Goodness either way, there must exist a fixed value  $F^j[X_{-\ell}^{n_{A_j}}]$  such that Subconditions (2) and (3) are true for  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$ .

We now show that the value  $F_\ell^U[X_{-\ell}^{n_U}] := \min_j(F_\ell^j[X_{-\ell}^{n_{A_j}}])$  causes Subconditions (2) and (3) to be true for the projection  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ , thus proving that this projection satisfies 1-Goodness.

**To show:**  $x_\ell < F_\ell^U[X_{-\ell}^{n_U}]$  implies  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell) = F_\ell^U[X_{-\ell}^{n_U}]$ . The condition  $x_\ell < F_\ell^U[X_{-\ell}^{n_U}]$  implies that for all  $j$ ,  $x_\ell < F_\ell^j[X_{-\ell}^{n_{A_j}}]$ . Then, for all  $j$ ,  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell) = F_\ell^j[X_{-\ell}^{n_{A_j}}]$  by Subcondition (2) for the various  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$ . Therefore,  $F_\ell^U[X_{-\ell}^{n_U}] = \min_j(F_\ell^j[X_{-\ell}^{n_{A_j}}]) = \min_j(T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)) = T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ , where the last step is by Eqn (4). This establishes Subcondition (2) for the projection  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ .

**To show:**  $x_\ell \geq F_\ell^U[X_{-\ell}^{n_U}]$  implies  $x_\ell \geq T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ . Because  $x_\ell \geq F_\ell^U[X_{-\ell}^{n_U}] = \min_j(F_\ell^j[X_{-\ell}^{n_{A_j}}])$ , there exists a  $j$  such that  $x_\ell \geq F_\ell^j[X_{-\ell}^{n_{A_j}}]$ . By Subcondition (3) for this  $T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$ , we have  $x_\ell \geq T_\ell^{(j)}[X_{-\ell}^{n_{A_j}}](x_\ell)$ . By Eqn (4), we then have  $x_\ell \geq T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ , thus establishing Subcondition (3) for  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$ .

Finally, because the above argument applies to every projection  $T_\ell^U[X_{-\ell}^{n_U}](x_\ell)$  of  $T^U(X^{n_U})$ , we have proved the desired result that  $T^U(X^{n_U})$  satisfies condition 3.5.  $\square$

### 3.6 Unbiasedness of EstimateOnSubPopulation()

We now show that 1-Goodness of a TCF implies that the corresponding instantiation of the Theta-Sketch Framework provides unbiased estimates of the number of unique identifiers on a stream or on the union of multiple streams.

**Theorem 3.10.** *Let  $A$  be a stream containing  $n_A$  unique identifiers, and let  $P$  be a property evaluating to 1 on an arbitrary subset of the identifiers. Let  $h$  denote a random hash function. Let  $T$  be a threshold choosing function that satisfies Condition 3.5. Let  $(\theta, S_A)$  denote a sketch of  $A$  created by  $\text{samp}[T](k, A, h)$ , and as usual let  $P(S_A)$  denote the subset of hash values in  $S_A$  whose corresponding identifiers satisfy  $P$ . Then  $E_h(\hat{n}_{P,A}) := E_h\left(\frac{|P(S_A)|}{\theta}\right) = n_{P,A}$ .*

Theorems 3.9 and 3.10 together imply that, in the multi-stream setting, the estimate  $\hat{n}_{P,U}$  for  $n_{P,U}$  output by the Theta-Sketch Framework is unbiased, assuming the base sampling schemes  $\text{samp}_j()$  each use a TCF  $T^{(j)}$  satisfying 1-Goodness.

*Proof.* Let  $A$  be a stream, and let  $T$  be a Threshold Choosing Function that satisfies 1-Goodness. Fix any  $\ell \in A$ . For any assignment  $X^{n_A}$  of hash values to identifiers in  $A$ , define the ‘‘per-identifier estimate’’  $V_\ell$  as follows:

$$V_\ell(X^{n_A}) = \frac{S_\ell(X^{n_A})}{T(X^{n_A})} \quad \text{where} \quad S_\ell(X^{n_A}) = \begin{cases} 1 & \text{if } x_\ell < T(X^{n_A}) \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

Because  $T$  satisfies 1-Goodness, there exists a fixed threshold  $F(X_{-\ell}^{n_A})$  for which it is a straightforward exercise to verify that:

$$V_\ell(X^{n_A}) = \begin{cases} 1/F(X_{-\ell}^{n_A}) & \text{if } x_\ell < F(X_{-\ell}^{n_A}) \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Now, conditioning on  $X_{-\ell}^{n_A}$  and taking the expectation with respect to  $x_\ell$ :

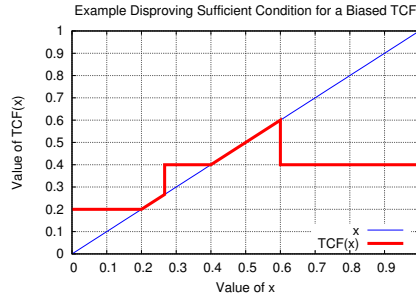
$$E(V_\ell|X_{-\ell}^{n_A}) = \int_0^1 V_\ell[X^{n_A}](x_\ell) dx_\ell = F(X_{-\ell}^{n_A}) \cdot \frac{1}{F(X_{-\ell}^{n_A})} = 1. \quad (7)$$

Since Equation (7) establishes that  $E(V_\ell) = 1$  when conditioned on each  $X_{-\ell}^{n_A}$ , we also have  $E(V_\ell) = 1$  when the expectation is taken over all  $X^{n_A}$ . By linearity of expectation, we conclude that  $E(\hat{n}_{P,A}) = \sum_{\ell \in A: P(\ell)=1} E(V_\ell) = n_{P,A}$ .  $\square$

**Is 1-Goodness Necessary for Unbiasedness?** Here we give an example showing that 1-Goodness cannot be substantially weakened while still guaranteeing unbiasedness of the estimate  $\hat{n}_{P,U}$  returned by the Theta-Sketch Framework. By construction, the following threshold choosing function causes the estimator of the Theta-Sketch Framework to be biased upwards.

$$T(X^{n_A}) = \begin{cases} m_k & \text{if } \frac{k-1}{m_k} > \frac{k}{m_{k+1}} \\ m_{k+1} & \text{otherwise} \end{cases} \quad (8)$$

Therefore, by the contrapositive of Theorem 3.10, it cannot satisfy Condition 3.5. It is an interesting exercise to try to establish this fact directly. It can be done by exhibiting a specific target size  $k$ , stream  $A$ , and partial assignment of hash values  $X_{-\ell}^{n_A}$  such that no fixed threshold  $F_\ell[X_{-\ell}^{n_A}]$  exists that would satisfy (2) and (3). Here is one such example:  $k = 3$ ,  $h(A) = \{0.1, 0.2, 0.4, 0.7, x_\ell\}$ .



The non-existence of the required fixed threshold is proved by the above plot of  $T(x_\ell)$ . The only value of  $F_\ell[X_{-\ell}^{n_A}]$  that would satisfy subcondition (2) is 0.2. However, that value does *not* satisfy (3), because  $T(x_\ell) > x_\ell$  for  $8/30 < x_\ell < 0.4$ .

### 3.7 1-Goodness and Monotonicity Imply Variance Bound

As usual, let  $U = \cup_{i=1}^m A_i$  be the union of  $m$  data streams. Our goal in this section is to identify conditions on a threshold choosing function which guarantee the following: whenever the Theta-Sketch Framework is instantiated with a TCF  $T$  satisfying the conditions, then for any property  $P \subseteq [n]$ , the variance  $\sigma^2(\hat{n}_{P,U})$  of the estimator obtained from the Theta-Sketch Framework is bounded above by the variance of the estimator obtained by running  $\text{samp}[T]()$  on the stream  $A^* := A_1 \circ A_2 \circ \dots \circ A_m$  obtained by concatenating  $A_1, \dots, A_m$ .

It is easy to see that 1-Goodness alone is not sufficient to ensure such a variance bound. Consider, for example, a TCF  $T$  that runs KMV on a stream  $A$  unless it determines that  $n_A \geq C$ , for some fixed value  $C$ , at which point it sets  $\theta$  to 1 (thereby causing  $\text{samp}[T]()$  to sample all elements from  $A$ ). Note that such a base sampling algorithm is not implementable by a sublinear space streaming algorithm, but  $T$  nonetheless satisfies 1-Goodness. It is easy to see that such a base sampling algorithm will fail to satisfy our desired comparative variance result when run on constituent streams  $A_1, \dots, A_m$  satisfying  $n_{A_i} < C$  for all  $i$ , and  $n_U > C$ . In this case, the variance of  $\hat{n}_U$  will be positive, while the variance of the estimator obtained by running  $\text{samp}[T]$  directly on  $A^*$  will be 0.

Thus, for our comparative variance result to hold, we assume that  $T$  satisfies both 1-Goodness and the following additional monotonicity condition.

**Condition 3.11 (Monotonicity Condition).** Let  $A_0, A_1, A_2$  be any three streams, and let  $A^* := A_0 \circ A_1 \circ A_2$  denote their concatenation. Fix any hash function  $h$  and parameter  $k$ . Let  $\theta = T(k, A_1, h)$ , and  $\theta' = T(k, A^*, h)$ . Then  $\theta' \leq \theta$ .

**Theorem 3.12.** Suppose that the Theta-Sketch Framework is instantiated with a TCF  $T$  that satisfies Condition 3.5 (1-Goodness), as well as Condition 3.11 (monotonicity). Fix a property  $P$ , and let  $A_1, \dots, A_m$ , be  $m$  input streams. Let  $U = \cup A_j$  denote the union of the distinct labels in the input streams. Let  $A^* = A_1 \circ A_2 \circ \dots \circ A_m$  denote the

concatenation of the input streams. Let  $(\theta^*, S^*) = \text{samp}[T](k, A^*, h)$ , and let  $\hat{n}_{P, A^*}^{A^*}$  denote the estimate of  $n_{P, A^*} = n_{P, U}$  obtained by evaluating  $\text{EstimateOnSubPopulation}((\theta^*, S^*), P)$ . Let  $(\theta^U, S^U) = \text{ThetaUnion}(\{(\theta_j, S_j)\})$ , and let  $\hat{n}_{P, U}^U$  denote the estimate of  $n_{P, U} = n_{P, A^*}$  obtained by evaluating  $\text{EstimateOnSubPopulation}((\theta^U, S^U), P)$ . Then, with the randomness being over the choice of hash function  $h$ ,  $\sigma^2(\hat{n}_{P, U}^U) \leq \sigma^2(\hat{n}_{P, A^*}^{A^*})$ .

The proof of Theorem 3.12 is somewhat involved, and is deferred to Appendix A.

**On the applicability of Theorem 3.12.** It is easy to see that Condition 3.11 holds for any TCF that is (1) order-insensitive and (2) has the property that adding another distinct item to the stream cannot increase the resulting threshold  $\theta$ . The TCF  $T$  used in multiKMV (namely,  $T(k, A, h) = m_{k+1}$ ), satisfies these properties, as does the TCF used in Adaptive Sampling. Since we already showed that both of these TCF's satisfy 1-Goodness, Theorem 3.12 applies to multiKMV and multiAdapt. In Section 3.9, we introduce the pKMV algorithm, which is useful in multi-stream settings where the distribution of stream lengths is highly skewed, and we show that Theorem 3.12 applies to this algorithm as well.

In Section 4, we introduce the Alpha Algorithm and show that it satisfies 1-Goodness. Unfortunately, the Alpha Algorithm does not satisfy monotonicity in general. The algorithm does, however, satisfy monotonicity under the promise that  $A_1, \dots, A_m$  are pairwise disjoint, and Theorem 3.12 applies in this case. Our experiments (Section 5.2) suggest that, in practice, the normalized variance in the multi-stream setting is not much larger than in the pairwise disjoint case.

### 3.8 Handling Set Intersections

The Theta-Sketch Framework can be tweaked in a natural way to handle set intersection and other set operations, just as was the case for multiKMV (cf. Section 6.2). Specifically, define  $\theta_U = \min_{j=1}^m \theta_j$ , and  $S_I = \{(x \in \cap_j S_j) < \theta_U\}$ . The estimator for  $n_{P, I}$  is  $\hat{n}_{P, I} := |P(S_I)|/\theta_U$ .

It is not difficult to see that  $\hat{n}_{P, I}$  is exactly equal to  $\hat{n}_{P', U}$ , where  $P'$  is the property that evaluates to 1 on an identifier if and only if the identifier satisfies  $P$  and is also in  $I$ . Since the latter estimator was already shown to be unbiased with variance bounded as per Theorem 3.12,  $\hat{n}_{P, I}$  satisfies the same properties.

### 3.9 The pKMV Variant of KMV

**Motivation.** An internet company involved in online advertising typically faces some version of the following problem: there is a huge stream of events representing visits of users to web pages, and a huge number of relevant “profiles”, each defined by the combination of a predicate on users and a predicate on web pages. On behalf of advertisers, the internet company must keep track of the count of distinct users who generate events that match each profile. The distribution (over profiles) of these counts typically is highly skewed and covers a huge dynamic range, from hundreds of millions down to just a few.

Because the summed cardinalities of all profiles is huge, the brute force technique (of maintaining, for each profile, a hash table of distinct user ids) would use an impractical amount of space. A more sophisticated approach would be to run multiKMV, treating each profile as separate stream  $A_i$ . This effectively replaces each hash table in the brute force approach with a KMV sketch. The problem with multiKMV in this setting is that, while KMV does avoid storing the entire data stream for streams containing more than  $k$  distinct identifiers, KMV produces no space savings for streams shorter than  $k$ . Because the vast majority of profiles contain only a few users, replacing the hash tables in the brute force approach by KMV sketches might still use an impractical amount of space.

On the other hand, fixed-threshold sampling with  $\theta = p$  for a suitable sampling rate  $p$ , would always result in an expected factor  $1/p$  saving in space, relative to storing the entire input stream. However, this method may result in too large a sample rate for long streams (i.e., for profiles satisfied by many users), also resulting in an impractical amount of space.

**The pKMV algorithm.** In this scenario, the hybrid Threshold Choosing Function  $T(k, A, h) = \min(m_{k+1}, p)$  can be a useful compromise, as it ensures that even short streams get downsampled by a factor of  $p$ , while long streams produce at most  $k$  samples. While it is possible to prove that this TCF satisfies 1-Goodness via a direct case analysis, the property can also be established by an easier argument: Consider a hypothetical computation in which the ThetaUnion

procedure is used to combine two sketches of the same input stream: one constructed by KMV with parameter  $k$ , and one constructed by fixed-threshold sampling with parameter  $p$ . Clearly, this computation outputs  $\theta = \min(m_{k+1}, p)$ . Also, since KMV and fixed-threshold sampling both satisfy 1-Goodness, and ThetaUnion preserves 1-Goodness (cf. Theorem 3.10),  $T$  also satisfies 1-Goodness.

It is easy to see that Condition 3.11 applies to  $T(k, A, h) = \min(m_{k+1}, p)$  as well. Indeed,  $T$  is clearly order-insensitive, so it suffices to show that adding an additional identifier to the stream cannot increase the resulting threshold. Since  $p$  never changes, the only way that adding another distinct item to the stream could increase the threshold would be by increasing  $m_{k+1}$ . However, that cannot happen.

## 4 Alpha Algorithm

### 4.1 Motivation and Comparison to Prior Art

Section 3’s theoretical results are strong because they cover such a wide class of base sampling algorithms. In fact, 1-Goodness even covers base algorithms that lack certain traditional properties such as invariance to permutations of the input, and uniform random sampling of the input. We are now going to take advantage of these strong theoretical results for the Theta Sketch Framework by devising a novel base sampling algorithm that lacks those traditional properties, but still satisfies 1-Goodness. Our main purpose for describing our Alpha Algorithm in detail is to exhibit the generality of the Theta-Sketch Framework. Nonetheless the Alpha Algorithm does have the following advantages relative to HLL, KMV, and Adaptive Sampling.

**Advantages over HLL.** Unlike HLL, the Alpha Algorithm provides unbiased estimates for  $\text{DISTINCT}_P$  queries for non-trivial predicates  $P$ . Also, when instantiating the Theta-Sketch Framework via the Alpha Algorithm in the multi-stream setting, the error behavior scales better than HLL for general set operations (cf. Section 2.2). Finally, because the Alpha Algorithm computes a sample, its output is human-interpretable and amenable to post-processing.

**Advantages over KMV.** Implementations of KMV must either use a heap data structure or quickselect [22] to give quick access to the  $k+1^{\text{st}}$  smallest unique hash value seen so far. The heap-based implementation yields  $O(\log k)$  update time, and quickselect, while achieving  $O(1)$  update time, hides a large constant factor in the Big-Oh notation (cf. Section 2.2). The Alpha Algorithm avoids the need for a heap or quickselect, yielding superior practical performance.

**Advantages over Adaptive Sampling.** The accuracy of Adaptive Sampling oscillates as  $n_A$  increases. The Alpha Algorithm avoids this behavior.

The remainder of this section provides a detailed analysis of the Alpha Algorithm. In particular, we show that it satisfies 1-Goodness, and we give quantitative bounds on its variance in the single-stream setting. Later (see Section 5.1), we describe experiments showing that, in both the single- and multi-stream settings, the Alpha Algorithm achieves a novel tradeoff between accuracy, space usage, update speed, and applicability.

**Detailed Section Roadmap.** Section 4.2 describes the threshold choosing function AlphaTCF that creates the instantiation of the Theta Sketch Framework whose base algorithm we refer to as the Alpha Algorithm. Section 4.3 establishes that AlphaTCF satisfies 1-Goodness, implying, via Theorem 3.10 that EstimateOnSubPopulation() is unbiased on single streams and on unions and intersections of streams in the framework instantiation created by plugging in AlphaTCF. Section 4.4 bounds the space usage of the Alpha Algorithm, as well as its variance in the single-stream setting. Section 4.5 discusses the algorithm’s variance in the multistream setting. Finally, Section 4.6 describes the HIP estimator derived from the Alpha Algorithm (see Section 6 for an introduction to HIP estimators).

### 4.2 AlphaTCF

Algorithm 2 describes the threshold choosing function AlphaTCF. AlphaTCF can be viewed as a tightly interleaved combination of two different processes. One process uses the set  $D$  to remove duplicate items from the raw input stream; the other process uses a technique similar to Approximate Counting [25] to estimate the number of items in the de-duped stream created by the first process. In addition, the second process maintains and frequently reduces a

---

**Algorithm 2** The Alpha Algorithm’s Threshold Choosing Function
 

---

```

1: Function AlphaTCF (target size  $k$ , stream  $A$ , hash function  $h$ )
2:  $\alpha \leftarrow k/(k+1)$ .
3:  $\text{prefix}(h(A)) \leftarrow$  shortest prefix of  $h(A)$  containing exactly  $k$  unique hash values.
4:  $\text{suffix}(h(A)) \leftarrow$  the corresponding suffix.
5:  $D \leftarrow$  the set of unique hash values in  $\text{prefix}(h(A))$ .
6:  $i \leftarrow 0$ .
7: for all  $x \in \text{suffix}(h(A))$  do
8:   if  $x < \alpha^i$  then
9:     if  $x \notin D$  then
10:       $i \leftarrow i + 1$ .
11:       $D \leftarrow D \cup \{x\}$ .
12:     end if
13:   end if
14: end for
15: return  $\theta \leftarrow \alpha^i$ .

```

---

threshold  $\theta = \alpha^i$  that is used by the first process to identify hash values that *cannot* be members of  $S$ , and therefore don’t need to be placed in the de-duping set  $D$ , thus limiting the growth of that set.

If the set  $D$  is implemented using a standard dynamically-resized hash table, then well-known results imply that the amortized cost<sup>6</sup> of processing each stream element is  $O(1)$ , and the space occupied by the hash table is  $O(|D|)$ , which grows logarithmically with  $n$ .

However, there is a simple optimized implementation of the Alpha Algorithm, based on Cuckoo Hashing, that implicitly, and at zero cost, deletes all members of  $D$  that are not less than  $\theta$ , and therefore are not members of  $S$  (see Section 5.1). This does not affect correctness, because those deleted members will not be needed for future de-duping tests of hash values that will all be less than  $\theta$ . Furthermore, in Theorem 4.2 below, it is proved that  $|S|$  is tightly concentrated around  $k$ . Hence, the space usage of this optimized implementation is  $O(k)$  with probability  $1 - o(1)$ .

### 4.3 AlphaTCF Satisfies 1-Goodness

We will now prove that AlphaTCF satisfies 1-Goodness.

**Theorem 4.1.** *If  $T(X^{n_A}) = \text{AlphaTCF}$ , then every fix-all-but-one projection  $T_\ell[X_{-\ell}^{n_A}](x_\ell)$  of  $T(X^{n_A})$  satisfies 1-Goodness.*

*Proof.* Fix the number of distinct identifiers  $n_A$  in  $A$ . Consider any identifier  $\ell$  appearing in the stream, and let  $x = h(\ell)$  be its hash value. Fix the hash values of all other elements of the sequence of values  $X_{-\ell}^{n_A}$ . We need to exhibit a threshold  $F$  such that  $x < F$  implies  $T_\ell[X_{-\ell}^{n_A}](x_\ell)(x) = F$  and  $x \geq F$  implies  $T_\ell[X_{-\ell}^{n_A}](x) \leq x$ .

First, if  $x$  lies in one of the first  $k + 1$  positions in the stream, then  $T_\ell[X_{-\ell}^{n_A}](x)$  is a constant independent of  $x$ ; in this case,  $F$  can be set to that constant.

Now for the main case, suppose that  $\ell$  does not lie in one of the first  $k + 1$  positions of the stream. Consider a subdivision of the hashed stream into the initial segment preceding  $x = h(\ell)$ , then  $x$  itself, then the final segment that follows  $x$ . Because all hash values besides  $x$  are fixed in  $X_{-\ell}^{n_A}$ , during the initial segment, there is a specific number  $a$  of times that  $\theta$  is decreased. When  $x$  is processed,  $\theta$  is decreased either zero or one times, depending on whether  $x < \alpha^a$ . Then, during the final segment,  $\theta$  will be decreased a certain number of additional times, where this number depends on whether  $x < \alpha^a$ . Let  $b$  denote the number of additional times  $\theta$  is decreased if  $x < \alpha^a$ , and  $c$  the number of additional times  $\theta$  is decreased otherwise. This analysis is summarized in the following table:

Rule	Condition on $x$	Final value of $\theta$
L	$x < \alpha^a$	$\alpha^{a+b+1}$
G	$x \geq \alpha^a$	$\alpha^{a+c+0}$

---

<sup>6</sup>Recent theoretical results imply that the update time can be made worst-case  $O(1)$  [2,3].

We prove the theorem using the threshold  $F = \alpha^{a+b+1}$ . We note that  $F = \alpha^{a+b+1} < \alpha^a$ , so  $F$  and  $\alpha^a$  divide the range of  $x$  into three disjoint intervals, creating three cases that need to be considered.

Case 1:  $x < F < \alpha^a$ . In this case, because  $x < F$ , we need to show that  $T_\ell[X_{-\ell}^{n_A}](x) = F$ . By Rule L,  $T_\ell[X_{-\ell}^{n_A}](x) = \alpha^{a+b+1} = F$ .

Case 2:  $F \leq x < \alpha^a$ . Because  $x \geq F$ , we need to show that  $T_\ell[X_{-\ell}^{n_A}](x) \leq x$ . By Rule L,  $T_\ell[X_{-\ell}^{n_A}](x) = \alpha^{a+b+1} = F \leq x$ .

Case 3:  $F < \alpha^a \leq x$ . Because  $x \geq F$ , we need to show that  $T_\ell[X_{-\ell}^{n_A}](x) \leq x$ . By Rule G,  $T_\ell[X_{-\ell}^{n_A}](x) = \alpha^{a+c+0} \leq \alpha^a \leq x$ .  $\square$

#### 4.4 Analysis of Alpha Algorithm on Single Streams

The following two theorems show that the Alpha Algorithm's space usage and single-stream estimation accuracy are quite similar to those of KMV. That means that it is safe to use the Alpha Algorithm as a drop-in replacement for KMV in a sketching-based big-data system, which then allows the system to benefit from the Alpha Algorithm's low update cost. See the Experiments in Section 5.1.

**Random Variables.** When Line 15 of Algorithm 2 is reached after processing a randomly hashed stream, the program variable  $i$  is governed by a random variable  $\mathcal{I}$ . Similarly, when Line 3 of Algorithm 1 is subsequently reached, the cardinality of the set  $S$  is governed by a random variable  $\mathcal{S}$ . The following two theorems characterize the distributions of  $\mathcal{S}$  and of the Theta Sketch Framework's estimator  $\mathcal{S}/(\alpha^{\mathcal{I}})$ . Specifically, Theorem 4.2 shows that the number of elements sampled by the Alpha Algorithm is tightly concentrated around  $k$ , and hence its space usage is concentrated around that of KMV. Theorem 4.3 shows that the variance of the estimate returned by the Alpha Algorithm is very close to that of KMV. Their proofs are rather involved, and are deferred to Appendices B.1 and B.2 respectively.

**Theorem 4.2.** *Let  $\mathcal{S}$  denote the cardinality of the set  $S$  computed by the Alpha Algorithm's Threshold Choosing Function (Algorithm 2). Then:*

$$\mathbb{E}(\mathcal{S}) = k. \quad (9)$$

$$\sigma^2(\mathcal{S}) < \frac{k}{2} + \frac{1}{4}. \quad (10)$$

**Theorem 4.3.** *Let  $\mathcal{S}$  denote the cardinality of the set  $S$  computed by the Alpha Algorithm's Threshold Choosing Function (Algorithm 2). Then:*

$$\sigma^2(\mathcal{S}/(\alpha^{\mathcal{I}})) = \frac{(2k+1)n_A^2 - (k^2+k)(2n_A-1) - n_A}{2k^2} \quad (11)$$

$$< \frac{n_A^2}{k - \frac{1}{2}}. \quad (12)$$

#### 4.5 Variance of the Alpha Algorithm in the Multi-Stream Setting

Unfortunately, the Alpha Algorithm does not satisfy monotonicity (Condition 3.11) in general, and hence Theorem 3.12 does not immediately imply variance bounds in the multi-stream setting. In fact, we have identified contrived examples in the multi-stream setting on which the variance of the Theta-Sketch Framework when instantiated with the TCF of the Alpha Algorithm is slightly larger than the hypothetical estimator obtained by running the Alpha Algorithm on the concatenated stream  $A_1 \circ \dots \circ A_m$  (the worst-case setting appears to be when  $A_1 \dots A_m$  are all permutations of each other).

However, we show in this section that the Alpha Algorithm does satisfy monotonicity under the promise that all constituent streams are pairwise disjoint. This implies the variance guarantees of Theorem 3.12 do apply to the Alpha Algorithm under the promise that  $A_1, \dots, A_m$  are pairwise disjoint. Our experiments in Section 5.2 suggest that, in practice, the normalized variance of the Alpha Algorithm in the multi-stream setting is not much larger than in the pairwise disjoint case.

**Theorem 4.4.** *The TCF computed by the Alpha Algorithm satisfies Condition 3.11 under the promise that the streams  $A_1, A_2, A_3$  appearing in Condition 3.11 are pairwise disjoint.*

*Proof.* Inspection of Algorithm 2 shows that the Alpha Algorithm never increases  $\theta$  while processing a stream. Therefore, processing  $A_3$  after  $A_2$  cannot increase  $\theta$  above the value that it had at the end of processing  $A_2$ . Hence, it will suffice to prove that  $T(A_2) \geq T(A_1 \circ A_2)$ . Referring to Line 15 of the pseudocode, we see that  $\theta = \alpha^I$ , where  $I$  is the final value of the program variable  $i$ , so it suffices to prove that  $I(A_2) \leq I(A_1 \circ A_2)$ .

We will compare two execution paths of the Alpha Algorithm. The first path results from processing  $A_2$  by itself. The second path results from processing  $A_1 \circ A_2$ . We will now index the sequence of hash values of  $h(A_1 \circ A_2)$  in a special way:  $x_0$  will be the first hash value that reaches Line 8 of the pseudocode during the first execution path (where  $A_2$  is processed by itself). Elements of  $h(A_1 \circ A_2)$  that follow  $x_0$  will be numbered  $x_1, x_2, \dots$ , while elements of  $h(A_1 \circ A_2)$  that precede  $x_0$  will be numbered  $\dots, x_{-2}, x_{-1}$ . We remark that the boundary between negative and positive indices does not coincide with the boundary between  $A_1$  and  $A_2$ .

For  $j \geq 0$ , let  $I(j)$  denote the value of the program variable  $i$  immediately before processing the hash value  $x_j$  on the first execution path ( $A_2$  alone), and let  $I'(j)$  denote the same quantity for the second execution path ( $A_1 \circ A_2$ ). We will prove by induction that for all  $j \geq 0$ ,  $I(j) \leq I'(j)$ . The base case is trivial: by construction of our indexing scheme, at position 0, execution path one has had no opportunities yet to increment  $i$ , while execution path two might have had some opportunities to increment  $i$ . Hence  $I(0) = 0$  while  $I'(0) \geq 0$ .

Now for the induction step. At position  $j$ ,  $I(j) \leq I'(j)$ , and the two values of  $i$  are both integers, so the only possible way for  $I(j+1) > I'(j+1)$  to occur would be for  $I(j) = I'(j)$ , and for the tests at Line 8 and Line 9 of the pseudocode to both *pass* on the first execution path, while at least one of them *fails* on the second execution path. However, the test in Line 8 must have the same outcome for both paths, since they are comparing the same hash value  $x_j$  against the same threshold  $\alpha^i = \alpha^{i'}$ . Also, given the assumption that  $A_1$  and  $A_2$  are disjoint, the ‘‘novelty test’’ in Line 9 is determined solely by novelty within  $A_2$ . Hence, it must have the same outcome on both paths. We conclude that it is impossible for  $i$  to be incremented on the first path but not on the second path, so  $I(j+1) > I'(j+1)$  is impossible.  $\square$

## 4.6 HIP estimator

For single streams, the HIP estimator (see Section 6 for an introduction to HIP estimators) derived from the Alpha Algorithm turns out to equal  $k/\alpha^i$ . This estimator does not involve the size of the sample set  $S$ , and is therefore not the same thing as the estimator  $|S|/\alpha^i$  derived by instantiating the Theta-Sketch Framework with the Alpha Algorithm. The following theorem shows that the variance bound of the HIP estimator guaranteed by Theorem 4.5 is smaller than the variance bound for the vanilla Alpha Algorithm (cf. Theorem 4.3) by a factor of 2. It can be proved by using the analysis of Approximate Counting of [12, 25], or by using the analysis of HIP estimators of [7, 28]. To keep the paper self-contained, Appendix B.4 contains a proof of this result that utilizes several immediate results developed in the proof of Theorem 4.3.

**Theorem 4.5.** *Let  $n_A$  denote the number of distinct elements in stream  $A$ . If  $\alpha^i = \text{AlphaTCF}(k, A, h)$ , then:*

$$\begin{aligned} \mathbb{E}(k/\alpha^i) &= n_A, \\ \sigma^2(k/\alpha^i) &= \frac{n_A^2 - 2n_A k + k^2 - n_A + k}{2k} < \frac{n_A^2}{2k}, \\ \text{S.E.}(k/\alpha^i) &< 0.708/\sqrt{k}. \end{aligned}$$

## 5 Experiments

### 5.1 Single-Stream Experiments Using Synthetic Data

In this section we describe experiments using synthetic data showing that implementations of KMV, Adaptive Sampling, and the Alpha Algorithm can provide different tradeoffs between time, space, and accuracy in the single-stream setting.

All three implementations take advantage of a version of cuckoo hashing that treats as empty all slots containing hash values that are not less than the current value of  $\theta$ .

The code for our streaming implementation of the Alpha Algorithm closely resembles the pseudocode presented as Algorithm 2. The de-duping set  $D$  is stored in a cuckoo hash table that uses the just-mentioned self-cleaning trick. Hence  $D$  is in fact *always* equal to  $S$ , with no extra work needed in the form of table rebuilds or explicit delete operations.

Our implementation of Adaptive Sampling uses the same self-cleaning hash tables, but has a different rule for reducing  $\theta$ : multiply by 1/2 each time  $|S|$  reaches a pre-specified limit. Again, no delete operations or table rebuilds are needed, but this program needs to scan the table after each reduction in  $\theta$  to discover the current size of  $|S|$ .

Finally, our implementation of KMV again uses the same self-cleaning hash tables, but it also uses a heap to keep track of the current value of  $\theta = m_{k+1}$ . Hence it either uses more space than the other two algorithms, or it suffers from a reduction in accuracy due to sharing the space budget between the hash table and the heap. Also, it is slower than the other two algorithms because it performs heap operations in addition to hash table operations.

Our experiments compare the speed and accuracy on single streams of these implementations of the three algorithms. Accuracy was evaluated using the metric  $\sqrt{\text{meanSquaredError}}/n_A$ , measured during the course of 1 million runs of each algorithm. We employed two different sets of experimental conditions.

First, we compare under “equal- $k$ ” conditions, in which all three algorithms aim for  $|S| = t/2$ , where  $t = 2^{16}$  denotes the size of the hash table. Adaptive Sampling is configured to oscillate between roughly  $|S| = (1/3)t$  and  $|S| = (2/3)t$ . We remark that KMV consumes more space than the other two algorithms under these conditions because of its heap.

Second, we compare under “equal-space” conditions reflective of a live streaming system that needs to limit the amount of memory consumed by each sketch data structure. Under these conditions, KMV is forced to devote half of its space budget to the heap, while both Adaptive Sampling and the Alpha Algorithm are free to employ parameters that cause their hash tables to run at occupancy levels well over 1/2. In detail, for KMV  $|S| = (2/5)t$ , for the Alpha Algorithm  $|S| = (4/5)t$ , while Adaptive Sampling oscillates between roughly  $|S| = (2/5)t$  and  $|S| = (4/5)t$ .

Experimental results are plotted in Figure 2. Two things are obvious. First, the heap-based implementation of KMV is much slower than the other two algorithms. Second, the error curves of Adaptive Sampling have a strongly oscillating shape that can be undesirable in practice.

Under the equal- $k$  conditions, the error curves of KMV and the Alpha Algorithm are so similar that they cannot be distinguished from each other in the plot. However, under the equal space conditions, the Alpha Algorithm’s ability to operate at a high, steady occupancy level (of the hash table) causes its error to be the lowest of the three algorithms. This high, steady occupancy level also causes the Alpha Algorithm to be slightly slower than Adaptive sampling under these conditions, even though the latter needs to re-scan the table periodically, while the Alpha Algorithm does not.

## 5.2 A Multi-Stream Experiment Using Real Data

As discussed in Section 3.7, Theorem 3.12’s comparative variance result does not apply to the Alpha Algorithm in general. However, we proved in Section 4.1 that Theorem 3.12 does apply to the Alpha Algorithm when the input streams are disjoint. In this section we present empirical evidence suggesting that the Alpha Algorithm “almost” satisfies the variance bound of Theorem 3.12 on real data. Recall that Theorem 3.12 asserted that  $\sigma^2(\hat{n}_{P,U}^U) \leq \sigma^2(\hat{n}_{P,A^*}^{A^*})$  when the estimates are computed using TCFs satisfying 1-Goodness and monotonicity. Simplifying notation, and switching from variance to relative error, we will exhibit a scatter plot comparing  $\text{RE}_U(A_1, A_2)$  versus  $\text{RE}_{A^*}(A_1, A_2)$ , for numerous pairs  $(A_1, A_2)$  of sets from a naturally occurring dataset, using the TCF defined by the Alpha Algorithm. This scatter plot will show that only a tiny fraction of the pairs violates the bound asserted in the theorem.

**WebScope “Groups” Dataset.** This experiment is based on ydata-ygroups-user-group-membership-graph-v1.0, a dataset that is available from the Yahoo Research Alliance Webscope program. It contains anonymized and downsampled membership lists for about 640000 Yahoo Groups, circa 2005. Because of the downsampling, there are only about 1 million members in all. We restricted our attention to the roughly 10000 groups whose membership lists contained between 201 and 5429 members. Hence there were about 50 million pairs of groups to consider. Recalling that the comparative variance theorem applies to the Alpha Algorithm under the promise that groups are disjoint, we trimmed this set of 50 million pairs down to 5000 pairs that seemed most likely to violate the theorem because they had the



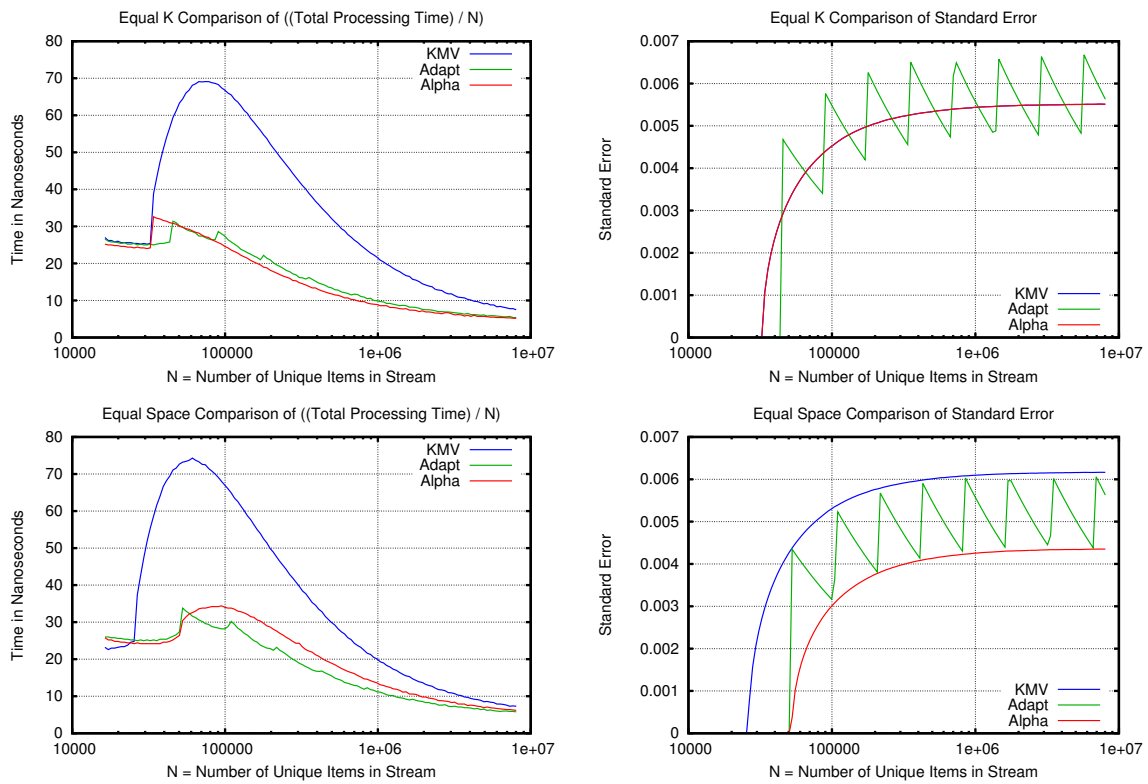


Figure 2: These plots illustrate the low stream processing cost and non-oscillating error curves of the Alpha Algorithm.

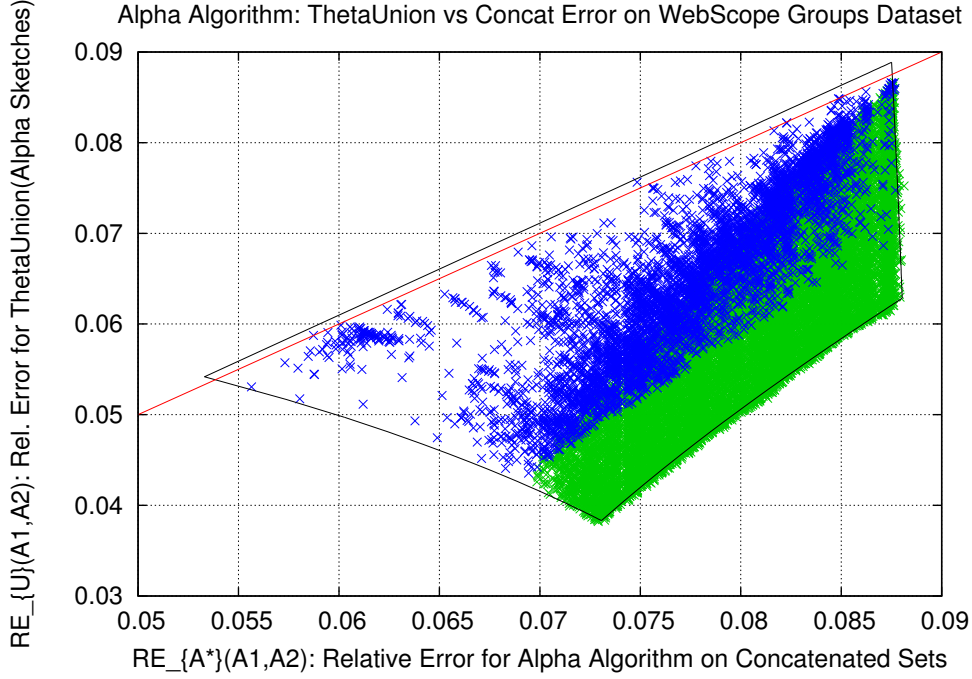


Figure 3: Most points are below the red line, showing that the comparative variance bound of Theorem 3.12 is “nearly true” for the Alpha Algorithm on the Webscope Dataset.

highest overlaps as measured by the similarity score  $\text{sim}(A_1, A_2) := (|A_1 \cap A_2| / \min(|A_1|, |A_2|))$ . We also examined another 13000 pairs of groups to fill out the scatter plot.

For each of these roughly 18000 pairs of groups, we empirically measured, by means of 100000 trials with  $k$  set to 128, the values of  $\text{RE}_U(A_1, A_2)$  and  $\text{RE}_{A^*}(A_1, A_2)$ , and plotted them in the scatter plot appearing as Figure 3. The 5000 high-overlap pairs are plotted in blue, while the other 13000 pairs are plotted in green. Strikingly, all but 2 of the roughly 18000 points lie on or below below the red line, thus indicating an outcome that is consistent with the comparative variance result. Because we included every pair of sets that had large overlap (as measured by  $\text{sim}(A_1, A_2)$ ) we conjecture that all of the other roughly 50 million pairs of sets also conform to the theorem.

Figure 3 also includes a heuristic “bounding box” plotted as a black quadrilateral. This bounding box was computed numerically from several ingredients. For the  $\text{RE}_U(A_1, A_2)$  side of the computation, we exploited the fact that for any given values of  $n$  and  $k$ , the Alpha Algorithm’s exact distribution over  $\theta$  values can be computed by dynamic programming using recurrences similar to the ones described in [12]. We also made the (counter-factual) assumption that three different hash functions are used to process the input sets  $A_1$  and  $A_2$ , and the output set  $S_U$ . This breaks the dependencies which complicate the analysis of the actual multi-stream instantiation of the Alpha Algorithm, in which only a single hash function is used during any given run. However, this counter-factual assumption also means that the resulting bounding box is not quite accurate. Finally, we did a grid search over all possible set-size pairs  $201 \leq |A_1| \leq |A_2| \leq 5429$  and all possible amounts of overlap, and traced out the boundary of the resulting combinations of computed relative errors. This boundary (see Figure 3) suggests that the comparative variance theorem is true for nearly all *possible* triples  $(|A_1|, |A_2|, |A_1 \cap A_2|)$  where  $201 \leq |A_1| \leq |A_2| \leq 5429$  and  $|A_1 \cap A_2| \leq |A_1|$ . Moreover, in those relatively few cases where the theorem is violated, the magnitude of the violation is small.

## 6 Detailed Overview of Prior Work

### 6.1 Algorithms for Single Streams

**HLL: HyperLogLog Sketches.** HLL is a sketching algorithm for the vanilla DISTINCT problem. It uses a hash function to randomly distribute the elements of a stream  $A$  amongst  $k$  buckets. For each bucket  $i$ , there is a register  $b_i$ , whose length is  $O(\log \log n)$  bits, that essentially contains the largest number of leading zeros in the hashed value of any stream element sent to that bucket. For each stream element, this data structure can clearly be updated in  $O(1)$  time. The HLL estimator for  $n_A$  which we denote  $\text{HLL}_A$ , is a certain non-linear function of the  $k$  bucket values  $b_i$ ; see [14]. It has been proved by [14] that, as  $n_A \rightarrow \infty$ ,  $E(\text{HLL}_A) \rightarrow n_A$ , and  $\sigma^2(\text{HLL}_A) \rightarrow 1.04(n_A^2/k)$ .

Unlike the KMV and Adaptive Sampling algorithms described below, it is not known how to extend the HLL sketch to estimate  $n_{P,A}$  for general properties  $P$  (unless, of course,  $P$  is known prior to stream processing). Qualitatively, the reason that HLL cannot estimate  $n_{P,A}$  is that, unlike the other algorithms, HLL does not maintain any kind of sample of identifiers from the stream.

**KMV: K'th Minimum Value Sketches.** The KMV sketching procedure for estimating  $\text{DISTINCT}(A)$  works as follows. While processing an input stream  $A$ , KMV keeps track of the set  $S$  of the  $k$  smallest unique hashed values of stream elements. The update time of a heap-based implementation of KMV is  $O(\log k)$ . The KMV estimator for  $\text{DISTINCT}(A)$  is

$$\text{KMV}_A = k/m_{k+1}, \quad (13)$$

where  $m_k$  denotes the  $k$ 'th smallest hash value. It has been proved by [5], [19], and others, that  $E(\text{KMV}_A) = n_A$ , and

$$\sigma^2(\text{KMV}_A) = \frac{n_A^2 - k n_A}{k - 1} < \frac{n_A^2}{k - 1}. \quad (14)$$

Duffield et al. [11] proposed to change the heap-based implementation of priority sampling to an implementation based on quickselect [22]. The same idea applies to KMV, which is a special case of priority sampling, and it reduces the sketch update cost from  $O(\log k)$  to amortized  $O(1)$ . However, this  $O(1)$  has a larger constant factor than that of competing methods.

The KMV sketching procedure can be extended to estimate  $n_{P,A}$  for any property  $P \subseteq [n]$ , as explained below. To accomplish this, the KMV sketch must keep not just the  $k$  smallest unique hash values that have been observed in the stream, but also the actual item identifiers corresponding to the hash values.<sup>7</sup> This allows the algorithm to determine which of the items in the sample satisfy the property  $P$ , even when  $P$  is not known until query time.

Motivated by the identity  $n_{P,A} = n_A \cdot (n_{P,A}/n_A)$ , the quantity  $\text{KMV}_A \cdot \text{est}(n_{P,A}/n_A)$  is a plausible estimate of  $n_{P,A}$ , for any sufficiently accurate estimate  $\text{est}(n_{P,A}/n_A)$  of  $n_{P,A}/n_A$ . Let  $S_A$  denote the  $k$  smallest unique hashed values in  $A$ , and recall (cf. Section 2.1) that  $P(S_A)$  denotes the subset of hash values in  $S_A$  whose corresponding identifiers in  $[n]$  satisfy the predicate  $P$  (the reason we require the sketch to store the actual identifiers that hashed to each value is to allow  $S_A$  to be determined from the sketch). Then the fraction  $|P(S_A)|/|S_A|$  can serve as the desired estimate of the fraction  $n_{P,A}/n_A$ . Essentially because  $S_A$  is a uniform random sample of  $A$ , it can be proved that the estimate  $\text{KMV}_{P,A} = \text{KMV}_A \cdot |P(S_A)|/|S_A|$  of  $n_{P,A}$  is unbiased, and has the following variance:<sup>8</sup>

$$\sigma^2(\text{KMV}_{P,A}) = \frac{n_{P,A}(n_A - k)}{k - 1} < \frac{n_{P,A} n_A}{k - 1}. \quad (15)$$

**Adaptive Sampling.** Adaptive Sampling maintains a sampling level  $i \geq 0$ , and the set  $S$  of all hash values less than  $2^{-i}$ ; whenever  $|S|$  exceeds a pre-specified size limit,  $i$  is incremented and  $S$  is scanned discarding any hash value that is now too big. Because a simple scan is cheaper than running quickselect, an implementation of this scheme can be cheaper than KMV. The estimator of  $n_A$  is  $\text{Adapt}_A = |S|/2^{-i}$ . It has been proved by [13] that this estimator is unbiased, and that  $\sigma^2(\text{Adapt}_A) \approx 1.44(n_A^2/(k - 1))$ , where the approximation sign hides oscillations caused by

<sup>7</sup>Technically, the sketch need not store the hash values if it stores the corresponding identifiers. Nonetheless, storing the hash values is often desirable in practice, to avoid the need to repeatedly evaluate the hash function.

<sup>8</sup>[5] analyzed the closely related estimator  $\text{KMV}'_{P,A} = \text{KMV}_A \cdot |P(S'_A)|/|S'_A|$ , where  $S'_A = S_A \cup \{m_{k+1}\}$ , proving unbiasedness and deriving the variance  $\sigma^2(\text{KMV}'_{P,A}) = (n_{P,A}((k + 1) n_A - (k + 1)^2 - n_A + k + 1 + n_{P,A})) / ((k + 1)(k - 1))$ .

the periodic culling of  $S$ . Like KMV, Adaptive Sampling can be extended to estimate  $n_{P,A}$  for any property  $P$ , via  $\text{Adapt}_{P,A} = \text{Adapt}_A \cdot |P(S_A)|/|S_A|$ . Note that, just as for KMV, this extension requires storing not just the hash values in  $S$ , but also the actual identifiers corresponding to each hash value.

Although the stream processing speed of Adaptive Sampling is excellent, the fact that its accuracy oscillates as  $n_A$  increases is a shortcoming of the method.

## 6.2 Algorithms for Set Operations on Multiple Streams

### HLL Sketches for Multiple Streams.

- **Set Union.** A sketch of  $U$  can be constructed from  $m$  HLL sketches of the  $A_j$ 's by taking the maximum of the  $m$  register values for each of the  $k$  buckets. The resulting sketch is identical to an HLL sketch constructed directly from  $U$ , so  $E(\text{HLL}_U) \rightarrow n_U$ , and  $\sigma^2(\text{HLL}_U) \rightarrow 1.04(n_U^2/k)$ .
- **Set Intersection.** Given constituent streams  $A_1, \dots, A_m$ , the HLL scheme can be extended via the Inclusion/Exclusion (IE) rule to estimate DISTINCT for various additional set-expressions other than set-union applied to  $A_1, \dots, A_m$ . This approach is awkward for complicated expressions, but is straightforward for simple expressions. For example, if  $m = 2$ , then the HLL+IE estimate of  $|I| = |A_1 \cap A_2|$  is  $\text{HLL}_{A_1} + \text{HLL}_{A_2} - \text{HLL}_U$ . Unfortunately, the variance of this estimate is approximately  $n_U^2/k$ . This is a factor of  $n_U^2/n_I^2$  larger than the variance of roughly  $n_I^2/k$  if one could somehow run HLL directly on  $I$ , and a factor of  $n_U/n_I$  worse than the variance achieved by the multiKMV algorithm described below. When  $n_I \ll n_U$ , this penalty factor overwhelms HLL's fundamentally good accuracy per bit.

In summary, the main limitations of HLL are its bad error scaling behavior when dealing with set operations other than set-union, as well as the inability to estimate  $\text{DISTINCT}_P$  queries for general properties  $P$ , even for a single stream  $A$ .

### multiKMV: KMV for Multiple Streams.

- **Set Union.** For any property  $P$ , there are two natural ways to extend KMV to estimate  $n_{P,U}$ , given a KMV sketch  $S_j$  containing the  $k + 1$  smallest unique hash values for each constituent stream  $A_j$ . The first is to use a “non-growing” union rule, and the second is to use a “growing” union rule (our term).

With the non-growing union rule, the sketch of  $U$  is simply defined to be the set of  $k + 1$  smallest unique hash values in  $\cup_{j=1}^m S_j$ . The resulting sketch is identical to a KMV sketch constructed directly from  $U$ , so  $E(\text{KMV}_U) = n_U$ , and  $\sigma^2(\text{KMV}_U) < n_U^2/(k - 1)$ . Just as the KMV sketch for a single stream  $A$  can be adapted to estimate  $n_{P,A}$  for any property  $P$ , this multi-stream variant of KMV can be adapted to provide an estimate  $\text{KMV}_{P,U}$  of  $n_{P,U}$ .

The growing union rule was introduced by Cohen and Kaplan [9]. This rule decreases the variance of estimates for unions and for other set expressions, but also increases the space cost of computing those estimates. Throughout, we refer to Cohen and Kaplan's algorithm as multiKMV. For each KMV input sketch  $S_j$ , let  $M_j$  denote that sketch's value of  $m_{k+1}$ . Define  $M_U = \min_{j=1}^m M_j$ , and  $S_U = \{x \in \cup_j S_j : x < M_U\}$ . Then  $n_U$  is estimated by  $\text{multiKMV}_U := |S_U|/M_U$ , and  $n_{P,U}$  is estimated by  $\text{multiKMV}_{P,U} := \text{multiKMV}_U \cdot |P(S_U)|/|S_U| = |P(S_U)|/M_U$ . [9] proved that  $\text{multiKMV}_{P,U}$  is unbiased and has variance that dominates the variance of the “non-growing” estimator  $\text{KMV}_{P,U}$ :

$$\sigma^2(\text{multiKMV}_{P,U}) \leq \sigma^2(\text{KMV}_{P,U}). \quad (16)$$

- **Set Intersection.** multiKMV can be tweaked in a natural way to handle set intersection and other set operations. Specifically, as in the set-union case, define  $M_U = \min M_j$ , and  $S_U = \{x \in \cup_j S_j : x < M_U\}$ . In addition, define  $S_I = \{(x \in \cap_j S_j) < M_U\}$ . The estimator for  $n_{P,I}$  is  $\text{multiKMV}_{P,I} := \text{multiKMV}_U \cdot |P(S_I)|/|S_U| = |P(S_I)|/M_U$ . It is not difficult to see that  $\text{multiKMV}_I$  is exactly equal to  $\text{multiKMV}_{P',U}$ , where  $P' = P \cap I$  is the property that evaluates to 1 on an identifier if and only if the identifier satisfies  $P$  and is also in  $I$ . Since the latter estimator was already shown to be unbiased with variance bounded as per Equation (1),  $\text{multiKMV}_{P,I}$  satisfies the same properties.

## multiAdapt: Adaptive Sampling for Multiple Streams.

- **Set Union.** Just as with KMV, for any property  $P$ , there are two natural ways to extend Adaptive Sampling to estimate  $n_{P,U}$ , given an Adaptive Sampling sketch  $S_j$  for each constituent stream  $A_j$ . The first is to use a non-growing union rule, and the second is to use a growing union rule. For brevity, we will only discuss the growing union rule, as proposed by [18]. We refer to this algorithm as multiAdapt. Let  $(i_j, S_j)$  be the sketch of the  $j$ 'th input stream  $A_j$ . The union sketch constructed from these sketches is  $(i_U = \max i_j, S_U = \{x \in \cup S_j : x < 2^{-i_U}\})$ . Then  $n_U$  is estimated by  $\text{multiAdapt}_U := |S_U|/2^{-i_U}$ , and  $n_{P,U}$  is estimated by  $\text{multiAdapt}_{P,U} := \text{multiAdapt}_U \cdot |P(S_U)|/|S_U|$ . [18] proved epsilon-delta bounds on the error of the estimator  $\text{multiAdapt}_{P,U}$ , but did not derive expressions for mean or variance. However, multiAdapt and multiKMV are in fact both special cases of our Theta-Sketch Framework, and in Section 3 of this paper we will prove (apparently for the first time) that  $\text{multiAdapt}_{P,U}$  is unbiased.
- **Set Intersection.** To our knowledge, prior work has not considered extending multiAdapt to handle set operations other than set-union on constituent streams. However, it is possible to tweak multiAdapt in a manner similar to multiKMV to handle these operations.

## 6.3 Other Related Work

Estimating the number of distinct values for data streams is a well studied problem. The problem of estimating result sizes of set expressions over multiple streams was concretely formulated by Ganguly et al. [16]. Motivated by the question of handling streams containing both insertions and deletions, their construction involves a 2-level hash function that essentially stores a set of counters for each bit-position of an HLL-type hash, and hence is inherently more resource intensive, both in terms of the space and update times.

$K$ 'th Minimum Value sketches were introduced by Bar-Yossef et al. [4], and developed into an unbiased scheme that handles set expressions by Beyer et al. [5]. Our own scheme is closely related to the schemes proposed and analyzed in Cohen and Kaplan [9], and in Gibbons and Tirthapura [18]. Chen, Cao and Bu [6] propose a somewhat different scheme for estimating unique counts with set expressions that is based on a data-structure related to the ‘‘probabilistic counting’’ sketches of [15], and also to the multi-bucket KMV sketches of [19] (with  $K = 1$ ). However, the guarantees proved by [6] are asymptotic in nature, and their system’s union sketches are the same size as base sketches, and therefore do not provide the increased accuracy that is possible with a ‘‘growing’’ union rule as in [9], in [18], and in this paper’s scheme.

Bottom-k sketches [8, 9] are a weighted generalization of KMV that provides unbiased estimates of the weights of arbitrary subpopulations of identifiers. They have small errors even under 2-independent hashing [27]. A closely related method for estimating subpopulation weights is priority sampling [11]. Although this paper’s Theta-Sketch Framework offers a broad generalization of KMV, it is not clear that it can support the entire generality of bottom-k sketches for weighted sets.

This paper’s ‘‘Alpha Algorithm’’ is inspired by the elegant *Approximate Counting* method of Morris [25], that has previously been applied to the estimation of the frequency moments  $F_p$ , for  $p \geq 1$ . By contrast, *our* task is to estimate  $\text{DISTINCT}_P$ . The Alpha Algorithm is able to do this because its Approximate Counting process is tightly interleaved with another process that removes duplicates from the input stream while maintaining a small memory footprint by using feedback from the approximate counter.

Kane et al. [23] gave a streaming algorithm for the  $\text{DISTINCTELEMENTS}$  problem that outputs a  $(1 + \epsilon)$ -approximation with constant probability, using  $\Theta(\epsilon^{-2} + \log(n))$  bits of space. This improves over the bit-complexity of HLL by roughly a  $\log \log n$  factor (and avoids the assumption of truly random hash functions). Like HLL, it is not known how to extend the algorithm to handle  $\text{DISTINCTONSUBPOPULATION}_P$  queries for non-trivial properties  $P$ , and the algorithm does not appear to have been implemented [21].

Tirthapura and Woodruff [29] give sketching algorithms for estimating  $\text{DISTINCTONSUBPOPULATION}_P$  queries for a special class of properties  $P$ . Specifically, they consider streams that contain tuples of the form  $(x, y)$ , where  $y$  is a numerical parameter, and the subpopulation  $P$  is specified via a lower or upper bound on  $y$ .

In very recent work, Cohen [7] and Ting [28] have proposed new estimators for  $\text{DISTINCTELEMENTS}$  (called ‘‘Historical Inverse Probability’’ (HIP) estimators in [7]). Any sketch which is generated by hashing of each element

in the data stream and is not affected by duplicate elements (such as HLL, KMV, Adaptive Sampling, and our Alpha Algorithm) has a corresponding HIP estimator, and [7, 28] show that the HIP estimator reduces the variance of the original sketching algorithm by a factor of 2. However, HIP estimators, in general, can only be computed when processing the stream, and this applies in particular to the HIP estimators of KMV and Adaptive Sampling. Hence, they do not satisfy the mergeability properties necessary to apply to multi-stream settings.

## References

- [1] N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.*, 58(1):137–147, 1999.
- [2] Y. Arbitman, M. Naor, and G. Segev. De-amortized cuckoo hashing: Provable worst-case performance and experimental results. In S. Albers, A. Marchetti-Spaccamela, Y. Matias, S. E. Nikolettseas, and W. Thomas, editors, *Automata, Languages and Programming, 36th International Colloquium, ICALP 2009, Rhodes, Greece, July 5-12, 2009, Proceedings, Part I*, volume 5555 of *Lecture Notes in Computer Science*, pages 107–118. Springer, 2009.
- [3] Y. Arbitman, M. Naor, and G. Segev. Backyard cuckoo hashing: Constant worst-case operations with a succinct representation. In *51th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2010, October 23-26, 2010, Las Vegas, Nevada, USA*, pages 787–796. IEEE Computer Society, 2010.
- [4] Z. Bar-Yossef, T. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream. In *Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer, 2002.
- [5] K. Beyer, R. Gemulla, P. J. Haas, B. Reinwald, and Y. Sismanis. Distinct-value synopses for multiset operations. *Communications of the ACM*, 52(10):87–95, 2009.
- [6] A. Chen, J. Cao, and T. Bu. A simple and efficient estimation method for stream expression cardinalities. In *Proceedings of the 33rd international conference on Very large data bases*, pages 171–182. VLDB Endowment, 2007.
- [7] E. Cohen. All-distances sketches, revisited: HIP estimators for massive graphs analysis. In R. Hull and M. Grohe, editors, *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS’14, Snowbird, UT, USA, June 22-27, 2014*, pages 88–99. ACM, 2014.
- [8] E. Cohen and H. Kaplan. Summarizing data using bottom-k sketches. In *Proceedings of the Twenty-sixth Annual ACM Symposium on Principles of Distributed Computing, PODC ’07*, pages 225–234, New York, NY, USA, 2007. ACM.
- [9] E. Cohen and H. Kaplan. Leveraging discarded samples for tighter estimation of multiple-set aggregates. In *Proceedings of the eleventh international joint conference on Measurement and modeling of computer systems*, pages 251–262. ACM, 2009.
- [10] G. Cormode. Sketch techniques for massive data. In G. Cormode, M. Garofalakis, P. Haas, and C. Jermaine, editors, *Synopses for Massive Data: Samples, Histograms, Wavelets and Sketches*, Foundations and Trends in Databases. NOW publishers, 2011.
- [11] N. G. Duffield, C. Lund, and M. Thorup. Priority sampling for estimation of arbitrary subset sums. *Journal of the ACM*, 54(6), 2007.
- [12] P. Flajolet. Approximate counting: a detailed analysis. *BIT Numerical Mathematics*, 25(1):113–134, 1985.
- [13] P. Flajolet. On adaptive sampling. *Computing*, 43(4):391–400, 1990.
- [14] P. Flajolet, É. Fusy, O. Gandouet, and F. Meunier. Hyperloglog: the analysis of a near-optimal cardinality estimation algorithm. *DMTCS Proceedings*, 0(1), 2008.

- [15] P. Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *Journal of computer and system sciences*, 31(2):182–209, 1985.
- [16] S. Ganguly, M. Garofalakis, and R. Rastogi. Processing set expressions over continuous update streams. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*, pages 265–276. ACM, 2003.
- [17] P. B. Gibbons. Distinct-values estimation over data streams. *Data Stream Management: Processing High-Speed Data Streams*, M. Garofalakis, J. Gehrke, and R. Rastogi, Eds. Springer, New York, NY, USA, 2007.
- [18] P. B. Gibbons and S. Tirthapura. Estimating simple functions on the union of data streams. In *Proceedings of the thirteenth annual ACM symposium on Parallel algorithms and architectures*, pages 281–291. ACM, 2001.
- [19] F. Giroire. Order statistics and estimating cardinalities of massive data sets. *Discrete Applied Mathematics*, 157(2):406–427, 2009.
- [20] A. Gronemeier and M. Sauerhoff. Applying approximate counting for computing the frequency moments of long data streams. *Theory of Computing Systems*, 44(3):332–348, 2009.
- [21] S. Heule, M. Nunkesser, and A. Hall. Hyperloglog in practice: Algorithmic engineering of a state of the art cardinality estimation algorithm. In *EDBT '13*, pages 683–692, New York, NY, USA, 2013. ACM.
- [22] C. A. R. Hoare. Algorithm 65: Find. *Communications of the ACM*, 4(7):321–322, July 1961.
- [23] D. M. Kane, J. Nelson, and D. P. Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the twenty-ninth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 41–52. ACM, 2010.
- [24] D. E. Knuth. *The Art of Computer Programming, Volume 3: (2Nd Ed.) Sorting and Searching*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 1998.
- [25] R. Morris. Counting large numbers of events in small registers. *Communications of the ACM*, 21(10):840–842, 1978.
- [26] S. Muthukrishnan. *Data streams: Algorithms and applications*. Now Publishers Inc, 2005.
- [27] M. Thorup. Bottom-k and priority sampling, set similarity and subset sums with minimal independence. In *Proceedings of the Forty-fifth Annual ACM Symposium on Theory of Computing, STOC '13*, pages 371–380, New York, NY, USA, 2013. ACM.
- [28] D. Ting. Streamed approximate counting of distinct elements: Beating optimal batch methods. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pages 442–451, New York, NY, USA, 2014. ACM.
- [29] S. Tirthapura and D. P. Woodruff. A General Method for Estimating Correlated Aggregates over a Data Stream. In *Proceedings of ICDE*, pages 162-173, 2012.

## A Proof of Theorem 3.12

### A.1 Proof Overview

The proof introduces the notion of the *fix-all-but-two projection* of a threshold choosing function  $T$ . We then introduce a new condition on TCF's that we call 2-Goodness (cf. Appendix A.2). On its face, 2-Goodness may appear to be a stronger requirement than 1-Goodness. However, we show in Section A.3 that this is not the case: 1-Goodness in fact

implies 2-Goodness.<sup>9</sup> We show in Appendix A.3 that 2-Goodness implies that “per-identifier estimates” output by the Theta-Sketch Framework are uncorrelated. Finally, in Section A.4, we use this result to complete the proof of Theorem 3.12.

## A.2 Definition of Fix-All-But-Two Projections and 2-Goodness

We begin by defining the Fix-All-But-Two Projection of a TCF.

**Definition A.1.** Let  $T$  be a threshold choosing function and fix a stream  $A$ . Let  $\ell_1 \neq \ell_2$  be two of the  $n_A$  unique identifiers in  $A$ . Let  $X_{-\ell_1, -\ell_2}^{n_A}$  be a fixed assignment of hash values to all unique identifiers in  $A$  except for  $\ell_1$  and  $\ell_2$ . Then the fix-all-but-two projection  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x_{\ell_1}, x_{\ell_2}) : [0, 1) \times [0, 1) \rightarrow (0, 1]$  of  $T$  is the function that maps values of  $(x_{\ell_1}, x_{\ell_2})$  to theta-sketch thresholds via the definition  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x_{\ell_1}, x_{\ell_2}) = T(X^{n_A})$ , where  $X^{n_A}$  is the obvious combination of  $X_{-\ell_1, -\ell_2}^{n_A}$ ,  $x_{\ell_1}$ , and  $x_{\ell_2}$ .

Next, we define the notion of 2-Goodness for bivariate functions.

**Definition A.2.** Let  $f(x, y) : [0, 1) \times [0, 1) \rightarrow (0, 1]$  be a bivariate function. We say that  $f$  satisfies 2-Goodness if there exists an  $F \in (0, 1]$  such that

- $\max(x, y) < F \Rightarrow f(x, y) = F$ .
- $\max(x, y) \geq F \Rightarrow f(x, y) \leq \max(x, y)$ .

Finally we are ready to define 2-Goodness for TCF’s.

**Condition A.3.** A threshold choosing function  $T(X^{n_A})$  satisfies 2-Goodness iff for every stream  $A$  containing  $n_A$  unique identifiers, every pair of identifiers  $\ell_1, \ell_2 \in A$ , and every fixed assignment  $X_{-\ell_1, -\ell_2}^{n_A}$  of hash values to the identifiers in  $A \setminus \{\ell_1, \ell_2\}$ , the fix-all-but-two projection  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x_{\ell_1}, x_{\ell_2})$  satisfies Definition A.2.

## A.3 1-Goodness Implies 2-Goodness

We are ready to show the (arguably surprising) result that if  $T$  satisfies 1-Goodness, then it also satisfies 2-Goodness.

**Theorem A.4.** Let  $T$  be a threshold choosing function that satisfies 1-Goodness. Then  $T$  also satisfies 2-Goodness.

*Proof.* Let  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}]$  be any fix-all-but-two projection of  $T$ . Notice that for any  $y' \in [0, 1)$ ,  $f(x) := T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y')$  is a fix-all-but-one projection of  $T$ . Similarly for any  $x' \in [0, 1)$ ,  $g(y) := T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x', y)$  is a fix-all-but-one-projection of  $T$ . Hence, 1-Goodness of  $T$  implies the following conditions hold:

**Property 1.** For each  $y' \in [0, 1)$ , there exists a  $G^{y'} \in (0, 1]$  such that:

- $x < G^{y'} \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y') = G^{y'}$ .
- $x \geq G^{y'} \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y') \leq x$ .

**Property 2.** For each  $x' \in [0, 1)$ , there exists a  $H^{x'} \in (0, 1]$  such that:

- $y < H^{x'} \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x', y) = H^{x'}$ .
- $y \geq H^{x'} \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x', y) \leq y$ .

To establish that  $T$  satisfies 2-Goodness, we want to prove that there exists an  $F \in (0, 1]$  such that

- $\max(x, y) < F \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = F$ .

<sup>9</sup>In fact, the two properties can be shown to be equivalent. We omit the reverse implication, since we will not require it to establish our variance bounds.



- $\max(x, y) \geq F \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) \leq \max(x, y)$ .

We will break the proof down into two lemmas.

**Lemma A.5.** *There exists an  $F \in (0, 1]$  such that  $\max(x, y) < F \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = F$ .*

*Proof.* By Property 1 above, there exists a  $G^0 \in (0, 1]$  such that

$$x < G^0 \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) = G^0. \quad (17)$$

Now consider any  $x$  in  $[0, G^0)$ . By Property 2 above, there exists a  $H^x \in (0, 1]$  such that:

$$y < H^x \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = H^x. \quad (18)$$

Plugging  $y = 0$  into Equation (18) gives  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) = H^x$ , while Equation (17) guarantees that  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) = G^0$ , so  $H^x = G^0$ . Substituting  $G^0$  into Equation (18) yields

$$y < G^0 \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = G^0. \quad (19)$$

Because  $x$  was any value in the interval  $[0, G^0)$ , the lemma is proved with  $F = G^0$ .  $\square$

**Lemma A.6.** *The threshold  $F$  whose existence was proved in Lemma A.5 also has the property that if  $\max(x, y) \geq F$ , then  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) \leq \max(x, y)$ .*

*Proof.* We start by assuming that  $\max(x, y) \geq F$ , so at least one of the following must be true:  $(x \geq F)$  or  $(y \geq F)$ . Without loss of generality we will assume that  $x \geq F$ . By Property 2 above, there exists an  $H^x \in (0, 1]$  such that

- $y < H^x \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = H^x$ .
- $y \geq H^x \Rightarrow T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) \leq y$ .

Our proof will have two cases, determined by whether  $y < H^x$  or  $y \geq H^x$ .

First case:  $y < H^x$ . In this case, because  $y < H^x$ ,  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = H^x$ . Also,  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) = H^x$ . But  $x \geq F = G^0$ , so  $T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) \leq x$ . Putting this all together gives:

$$T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) = H^x = T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, 0) \leq x \leq \max(x, y). \quad (20)$$

Second case:  $y \geq H^x$ . In this case, because  $y \geq H^x$ ,

$$T_{\ell_1, \ell_2}[X_{-\ell_1, -\ell_2}^{n_A}](x, y) \leq y \leq \max(x, y). \quad (21)$$

$\square$   
 $\square$

## 1-Goodness Implies Per-Identifier Estimates Are Uncorrelated

**Lemma A.7.** *Fix any stream  $A$ , threshold choosing function  $G$ , and pair  $\ell_1 \neq \ell_2$  in  $A$ . Define the “per-identifier estimates”  $V_{\ell_1}$  and  $V_{\ell_2}$  as in Equation (5). Then if  $T$  satisfies 1-Goodness, the covariance of  $V_{\ell_1}$  and  $V_{\ell_2}$  is 0. In symbols,*

$$\sigma(V_{\ell_1}, V_{\ell_2}) = E_{X^{n_A}}(V_{\ell_1} \cdot V_{\ell_2}) - E_{X^{n_A}}(V_{\ell_1}) \cdot E_{X^{n_A}}(V_{\ell_2}) = 0.$$

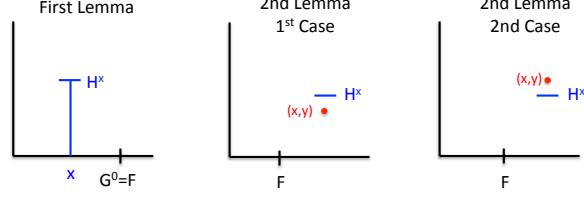


Figure 4: Some diagrams for Lemmas A.5 and A.6

*Proof.* Because  $T$  satisfies 1-Goodness, it also satisfies 2-Goodness (cf. Theorem A.4), and hence there exists a threshold  $F(X_{-\ell_1, -\ell_2}^{n_A})$  for which it is a straightforward exercise to verify that:

$$V_{\ell_1}(X^{n_A}) \cdot V_{\ell_2}(X^{n_A}) = \begin{cases} 1/F(X_{-\ell_1, -\ell_2}^{n_A})^2 & \text{if } \max(x_{\ell_1}, x_{\ell_2}) < F(X_{-\ell_1, -\ell_2}^{n_A}) \\ 0 & \text{otherwise.} \end{cases} \quad (22)$$

Now, conditioning on  $X_{-\ell_1, -\ell_2}^{n_A}$  and taking the expectation with respect to pairs  $(x_{\ell_1}, x_{\ell_2})$ :

$$E(V_{\ell_1} \cdot V_{\ell_2} | X_{-\ell_1, -\ell_2}^{n_A}) = \int_0^1 \int_0^1 V_{\ell_1}(X^{n_A}) V_{\ell_2}(X^{n_A}) dx_{\ell_1} dx_{\ell_2} = F(X_{-\ell_1, -\ell_2}^{n_A})^2 \cdot \frac{1}{F(X_{-\ell_1, -\ell_2}^{n_A})^2} = 1. \quad (23)$$

Since  $E(V_{\ell_1} V_{\ell_2} | X_{-\ell_1, -\ell_2}^{n_A}) = 1$  when conditioned on each  $X_{-\ell_1, -\ell_2}^{n_A}$ , we also have  $E(V_{\ell_1} V_{\ell_2}) = 1$  when the expectation is taken over all  $X^{n_A}$ . Meanwhile, since  $T$  satisfies 1-Goodness,  $E(V_{\ell_1}) = E(V_{\ell_2}) = 1$  (cf. Theorem 3.10). Hence,  $\sigma(V_{\ell_1}, V_{\ell_2}) = 0$ .  $\square$

As a corollary of Lemma A.7, we obtain the following result, establishing that the variance of  $\hat{n}_{P,A}$  is equal to the sum of the variances of the per-identifier estimates for all identifiers in  $A$  satisfying property  $P$ .

**Lemma A.8.** *Suppose that  $T$  satisfies 1-Goodness. Fix any stream  $A$ , and let  $\hat{n}_{P,A}$  denote the estimate for  $n_{P,A}$  obtained by running  $\text{samp}[T]()$  on  $A$  and feeding the resulting theta-sketch into  $\text{EstimateOnSubPopulation}()$ . Then*

$$\sigma^2(\hat{n}_{P,A}) = \sum_{\ell \in A: P(\ell)=1} \sigma^2(V_{\ell}).$$

*Proof.* Note that  $\hat{n}_{P,A} = \sum_{\ell \in A: P(\ell)=1} V_{\ell}$ . The claim then follows from Lemma A.7 combined with the fact that the variance of the sum of random variables equals the sum of the variances, provided that the variables appearing in the sum are uncorrelated.  $\square$

## A.4 Completing the Proof of Theorem 3.12

*Proof.* For every  $\ell$  that appears in the concatenated stream  $A^*$ , and for all  $X^{n_{A^*}}$ , we define the “per-identifier estimate”  $V_{\ell}(X^{n_{A^*}})$  as in Equation (5) with  $A = A^*$ , and relate it to the threshold  $F_{\ell}(X_{-\ell}^{n_{A^*}})$  as in Equation (6), also with  $A = A^*$ . It is then straightforward to verify that

$$\sigma^2(V_{\ell} | X_{-\ell}^{n_{A^*}}) = 1/F_{\ell}(X_{-\ell}^{n_{A^*}}) - 1. \quad (24)$$

Let  $T'$  be the TCF that was (implicitly) used to construct  $(\theta^U, S^U)$  from the  $m$  sketches of the individual streams  $A_j$ . By Theorem 3.9,  $T'$  satisfies 1-Goodness, so let  $F'_{\ell}(X_{-\ell}^{n_{A^*}})$  denote the corresponding threshold value for  $T'$  as in Equation (6). We claim that  $T'$  satisfies the following property:

$$\text{For all identifiers } \ell \in [n] \text{ and for all } X^{n_{A^*}}, F'_{\ell}(X_{-\ell}^{n_{A^*}}) \geq F_{\ell}(X_{-\ell}^{n_{A^*}}). \quad (25)$$

**Finishing the proof, assuming  $T'$  satisfies Property 25.** By Equation (24):

$$\sigma^2(V'_\ell | X_{-\ell}^{n_{A^*}}) \leq \sigma^2(V_\ell | X_{-\ell}^{n_{A^*}}). \quad (26)$$

Because this inequality holds for every specific  $X_{-\ell}^{n_{A^*}}$ , it also holds for any convex combination over  $X_{-\ell}^{n_{A^*}}$ 's, so

$$\sigma^2(V'_\ell) \leq \sigma^2(V_\ell).$$

Combining this with Lemma A.8, we conclude that

$$\sigma^2(\hat{n}_{P,U}^U) = \sum_{\ell \in A: P(\ell)=1} \sigma^2(V'_\ell) \leq \sum_{\ell \in A: P(\ell)=1} \sigma^2(V_\ell) = \sigma^2(\hat{n}_{P,A^*}^{A^*}).$$

**Proving that  $T'$  satisfies Property 25.** Fix any hash function  $h$ , which determines  $X^{n_U}$ , and also fixes hashed versions of the streams  $A_1, \dots, A_m$  and  $A^*$ . We will overload the symbols  $A_j$  and  $A^*$  to denote these hashed streams as well as the original streams. We need to prove that  $F'_\ell(X_{-\ell}^{n_U}) \geq F_\ell(A_{-\ell}^*)$ . This can be done in three steps. First, from the proof of Theorem 3.9 we know that there exists a  $j$  such that  $F'_\ell(X_{-\ell}^{n_U}) = F_\ell(A_{j,-\ell})$ . Second, because  $T$  satisfies 1-Goodness,  $F_\ell(A_{j,-\ell}) = T(Z(A_j, \ell))$  and  $F_\ell(A_{-\ell}^*) = T(Z(A^*, \ell))$ , where  $Z$  is a function that makes a copy of a hashed stream in which  $h(\ell)$  has been artificially set to zero. Third,  $Z(A^*, \ell)$  can be rewritten as the concatenation of 3 streams as follows:  $B_0 \circ Z(A_j, \ell) \circ B_2$ , where  $B_0 = Z(A_1, \ell) \circ Z(A_2, \ell) \circ \dots \circ Z(A_{j-1}, \ell)$ , and  $B_2 = Z(A_{j+1}, \ell) \circ \dots \circ Z(A_m, \ell)$ . Because  $T$  was assumed to satisfy the monotonicity condition, Condition 3.11, we then have

$$F'_\ell(X_{-\ell}^{n_U}) = T(Z(A_j, \ell)) \geq T(B_0 \circ Z(A_j, \ell) \circ B_2) = T(Z(A^*, \ell)) = F_\ell(A_{-\ell}^*). \quad (27)$$

□

## B Details of the Analysis of the Alpha Algorithm for Single Streams

### B.1 Proof of Theorem 4.2

Let  $S$  be the set produced by Line 3 of Algorithm 1 when AlphaTCF is plugged into the Theta Sketch Framework, and let  $\mathcal{S}$  be the random variable corresponding to  $|S|$ . In this section we compute  $E(\mathcal{S})$  and bound  $\sigma^2(\mathcal{S})$ .

We prove the top-level theorem using a lemma. The proofs of the theorem and lemma both involve two levels of conditioning. First we condition on the value  $\mathcal{I}$  of  $i$  when Line 15 of Algorithm 2 is reached. Then we further condition on  $\mathcal{J}^+$ , which we define to be the particular set of  $i$  stream positions on which increments occurred in Line 10 of Algorithm 2.

#### Restatement of Theorem 4.2.

$$E(\mathcal{S}) = k. \quad (28)$$

$$\sigma^2(\mathcal{S}) < \frac{k}{2} + \frac{1}{4}. \quad (29)$$

*Proof.* Using standard laws of probability, we perform the following decompositions:

$$E(\mathcal{S}) = \sum_i \Pr(\mathcal{I} = i) E(\mathcal{S} | \mathcal{I} = i) \quad (30)$$

$$E(\mathcal{S} | \mathcal{I} = i) = \sum_J \Pr(\mathcal{J}^+ = J | \mathcal{I} = i) E(\mathcal{S} | \mathcal{I} = i, \mathcal{J}^+ = J) \quad (31)$$

$$E(\mathcal{S}^2) = \sum_i \Pr(\mathcal{I} = i) E(\mathcal{S}^2 | \mathcal{I} = i) \quad (32)$$

$$E(\mathcal{S}^2 | \mathcal{I} = i) = \sum_J \Pr(\mathcal{J}^+ = J | \mathcal{I} = i) E(\mathcal{S}^2 | \mathcal{I} = i, \mathcal{J}^+ = J) \quad (33)$$

In Lemma B.1, we prove that for all  $i$  and  $J$ ,

$$E(\mathcal{S}|\mathcal{I} = i, \mathcal{J}^+ = J) = k. \quad (34)$$

Because this answer does not depend on  $J$ , the RHS of Equation (31) is a convex combination of equal values, so

$$E(\mathcal{S}|\mathcal{I} = i) = k.$$

Because this answer does not depend on  $i$ , the RHS of Equation (30) is a convex combination of equal values, so

$$E(\mathcal{S}) = k.$$

In Lemma B.1, we also prove that

$$E(\mathcal{S}^2|\mathcal{I} = i, \mathcal{J}^+ = J) = k^2 + \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2}.$$

Because this answer does not depend on  $J$ , the RHS of Equation (33) is a convex combination of equal values, so

$$E(\mathcal{S}^2|\mathcal{I} = i) = k^2 + \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2}. \quad (35)$$

This answer *does* depend on  $i$ , so we cannot use the exact same argument for a fourth time. However, with a little bit of algebra, one can go from Equation (35) to the inequality

$$E(\mathcal{S}^2|\mathcal{I} = i) < k^2 + \frac{k}{2} + \frac{1}{4},$$

whose RHS does not depend on  $i$ . Then the RHS of Equation (32) is a convex combination of values that are all less than  $k^2 + \frac{k}{2} + \frac{1}{4}$ . So:

$$\begin{aligned} E(\mathcal{S}^2) &< k^2 + \frac{k}{2} + \frac{1}{4}. \\ \sigma^2(\mathcal{S}) &< (k^2 + \frac{k}{2} + \frac{1}{4}) - k^2 = \frac{k}{2} + \frac{1}{4}. \end{aligned}$$

□

Now we will prove the lemma that was used above:

**Lemma B.1.** *For all  $i$  and  $J$ ,*

$$E(\mathcal{S}|\mathcal{I} = i, \mathcal{J}^+ = J) = k, \quad (36)$$

$$E(\mathcal{S}^2|\mathcal{I} = i, \mathcal{J}^+ = J) = k^2 + \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2}. \quad (37)$$

*Proof.* Because Line 9 of Algorithm 2 causes the algorithm to ignore duplicate labels, it will suffice to analyze a stream  $A$  of length  $n$  that doesn't contain any duplicates. Let  $h$  be a hash function that is chosen randomly. Let  $\{X_p | 1 \leq p \leq n\}$  be a sequence of  $n$  iid random variables, one per stream position, each drawn from the distribution Uniform(0,1). Let  $X^{n_A}$  be the cross product of the  $X_p$ 's; this random variable is our model of  $h(A)$ .  $\mathcal{I}$  is distributed as a random variable generated by first choosing a random  $X^{n_A}$ , then running Algorithm 2 on  $X^{n_A}$ , and then setting  $\mathcal{I}$  to be the value of the program variable  $i$  when Line 15 is reached. Define a set of  $n$  Bernoulli random variables  $S_p$ , one per stream position, derived from the variable  $\mathcal{I}$  and the variables  $X_p$  by the rule  $S_p = 1$  iff  $X_p < \alpha^i$ . Note that the  $S_p$ 's are *not* independent of each other. However, as we will see, they become independent after conditioning on the event  $(\mathcal{I} = i \text{ and } \mathcal{J}^+ = J)$ .

Now we will describe the effect of conditioning on both  $\mathcal{I} = i$  and  $\mathcal{J}^+ = J$ , by first describing how the original variables  $X_p$  are transformed into modified variables  $Y_p$  that are drawn from specific subintervals of  $(0, 1)$ . We will then introduce new Bernoulli variables  $S'_p$  defined by the rule  $S'_p = 1$  iff  $Y_p < \alpha^i$ , and finally compute the expected value and variance of  $(\mathcal{S}|\mathcal{I} = i, \mathcal{J}^+ = J) = \sum_p S'_p$ .

We are fixing a specific value  $i$  of  $\mathcal{I}$  and  $J$  of  $\mathcal{J}^+$ ; the latter is a size- $i$  subset of the set of  $n - k$  non-initial stream positions  $\{k + 1, k + 2, \dots, n - 1, n\}$ . The set of  $n - k - i$  non-initial stream positions that are not in  $J$  will be referred to as  $J^-$ . Let  $f(p, J)$  be the function that maps any non-initial position  $p$  to the number of non-initial positions before  $p$  that are members of  $J$ . We note that for  $p \in J$ ,  $f(p, J) \in \{0, 1, \dots, i - 1\}$ , and the mapping is one-to-one. For  $p \in J^-$ ,  $f(p, J) \in \{0, 1, \dots, i\}$ , and the mapping is not necessarily one-to-one.

Now we are ready to characterize the  $Y_p$ 's and the  $S'_p$ 's.

First let  $p$  be one of the  $k$  initial positions in the stream. In this case, conditioning on  $\mathcal{I} = i$  and  $\mathcal{J}^+ = J$  does not tell us anything about the value of  $X_p$ , so  $Y_p$  is drawn from the full interval  $(0, 1)$ , so  $\Pr(Y_p < \alpha^i) = \alpha^i$ ;  $E(S'_p) = \alpha^i$ , and  $\sigma^2(S'_p) = \alpha^i(1 - \alpha^i)$ .

Next, let  $p$  be one of the  $n - k - i$  positions in  $J^-$ . For this position, the test in Line 8 of Algorithm 2 failed, so we know that  $X_p \geq \alpha^{f(p, J)}$ , so  $Y_p$  is drawn uniformly from the interval  $[\alpha^{f(p, J)}, 1)$ . Because  $p \in J^-$ ,  $f(p, J) \leq i$ , so  $\alpha^i \leq \alpha^{f(p, J)} \leq Y_p$ , so  $\Pr(Y_p < \alpha^i) = 0$ ,  $E(S'_p) = 0$ , and  $\sigma^2(S'_p) = 0$ .

Finally, let  $p$  be one of the  $i$  positions that are in  $J$ . For this position, the test in Line 8 of Algorithm 2 succeeded, so we know that  $X_p < \alpha^{f(p, J)}$ , so  $Y_p$  is drawn uniformly from the interval  $(0, \alpha^{f(p, J)})$ , so  $\Pr(Y_p < \alpha^i) = \alpha^i / \alpha^{f(p, J)} = \alpha^{i-f(p, J)}$ . Now, because  $f(p, J)$  assumes each value in  $\{0, 1, \dots, i - 1\}$  as  $p$  is varied over the contents of  $J$ ,  $E(S'_p) = \Pr(Y_p < \alpha^i) = \alpha^{i-f(p, J)}$  assumes each value in  $\{\alpha^i, \alpha^{i-1}, \dots, \alpha^1\}$ . Similarly,  $\sigma^2(S'_p)$  assumes each value in  $\{\alpha^i(1 - \alpha^i), \dots, \alpha^1(1 - \alpha^1)\}$ . Note that the above analysis implies that the  $S'_p$ 's are independent of each other after conditioning on the event  $(\mathcal{I} = i$  and  $\mathcal{J}^+ = J)$ .

Putting together all of the above, and remembering that the random variables  $S'_p$  are independent due to the conditioning on  $\mathcal{I} = i$  and  $\mathcal{J}^+ = J$ :

$$E(\mathcal{S}|\mathcal{I} = i, \mathcal{J}^+ = J) = k \cdot \alpha^i + (n - k - i) \cdot 0 + \sum_{j=1}^i \alpha^j = k.$$

Here, we have used the fact that  $\alpha = \frac{k}{k+1}$ . In addition:

$$\begin{aligned} \sigma^2(\mathcal{S}|\mathcal{I} = i, \mathcal{J}^+ = J) &= k\alpha^i \cdot (1 - \alpha^i) + (n - k - i) \cdot 0 + \sum_{j=1}^i \alpha^j(1 - \alpha^j) \\ &= \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2}, \\ E(\mathcal{S}^2|\mathcal{I} = i, \mathcal{J}^+ = J) &= k^2 + \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2} \end{aligned}$$

□

## B.2 Proof of Theorem 4.3

**Preliminaries and Notation:** Because Line 9 of Algorithm 2 causes the algorithm to ignore duplicate labels, it will suffice to analyze streams that do not contain any duplicates.  $\mathcal{S}$  and  $\mathcal{I}$  are random variables giving the final values of  $|S|$  and  $i$  when Line 15 of Algorithm 2 is reached. Let  $\mathcal{Z} = \mathcal{S}/\alpha^{\mathcal{I}}$  denote the random variable for the estimate produced by the Theta Sketch framework when the Alpha Algorithm's TCF is used. It will be convenient to introduce a new variable  $u = n_A - k$  representing the number of stream items that are processed by the Alpha Algorithm *after*  $k$  initial items have been processed to initialize the set  $S$ . Recall that  $\alpha = k/(k + 1)$ .

**Restatement of Theorem 4.3.**

$$\sigma^2(\mathcal{Z}) = \frac{(2k + 1)n_A^2 - (k^2 + k)(2n_A - 1) - n_A}{2k^2} < \frac{n_A^2}{k - \frac{1}{2}}.$$

*Proof.* Due the 1-goodness of the Alpha Algorithm's TCF, we already know that  $\mathcal{Z}$  is an unbiased estimator for  $n_A$ , i.e., that  $E(\mathcal{Z}) = n_A = k + u$ . Hence,

$$\begin{aligned}\sigma^2(\mathcal{Z}) &= E(\mathcal{Z}^2) - E^2(\mathcal{Z}) \\ &= E(\mathcal{Z}^2) - (k + u)^2.\end{aligned}$$

Lemma B.2 (stated and proved below) gives a formula for  $E(\mathcal{Z}^2)$ . This allows us to complete the the analysis of  $\sigma^2(\mathcal{Z})$  as follows:

$$\begin{aligned}E(\mathcal{Z}^2) - (k + u)^2 &= \frac{k^2u + ku^2 + u(u-1)/2}{k^2} \\ &= \frac{(2k+1)n_A^2 - (k^2+k)(2n_A-1) - n_A}{2k^2} \\ &< \frac{n_A^2}{k - \frac{1}{2}}.\end{aligned}$$

□

**Lemma B.2.**

$$E(\mathcal{Z}^2) = \frac{k^2u + ku^2 + u(u-1)/2 + k^4 + 2k^3u + k^2u^2}{k^2}.$$

*Proof.*

$$\begin{aligned}E(\mathcal{Z}^2) &= \sum_{i=0}^u \sum_{s=0}^{k+i} \left(\frac{s}{\alpha^i}\right)^2 \Pr(\mathcal{S}=s|\mathcal{I}=i) \Pr(\mathcal{I}=i; u) \\ &= \sum_{i=0}^u \frac{1}{\alpha^{2i}} \Pr(\mathcal{I}=i; u) \sum_{s=0}^{k+i} s^2 \Pr(\mathcal{S}=s|\mathcal{I}=i) \\ &= \sum_{i=0}^u \frac{1}{\alpha^{2i}} \Pr(\mathcal{I}=i; u) E(\mathcal{S}^2|\mathcal{I}=i)\end{aligned}$$

Above, we use the somewhat onerous notation  $\Pr(\mathcal{I}=i; u)$  to emphasize that the distribution of  $\mathcal{I}$  depends on the fixed quantity  $u$  (i.e., on the number of distinct elements in the stream, minus  $k$ ). Making this dependence explicit will be useful later, when we analyze this distribution by establishing recurrences involving  $u$ .

A formula for  $E(\mathcal{S}^2|\mathcal{I}=i)$  appeared in Equation (35). Substituting this formula and continuing:

$$\begin{aligned}&= \sum_{i=0}^u \frac{1}{\alpha^{2i}} \Pr(\mathcal{I}=i; u) \left( \frac{\alpha - \alpha^{2i+1}}{1 - \alpha^2} + k^2 \right) \\ &= \frac{1}{1 - \alpha^2} \left[ \alpha \sum_{i=0}^u \frac{1}{\alpha^{2i}} \Pr(\mathcal{I}=i; u) - \sum_{i=0}^u \alpha \Pr(\mathcal{I}=i; u) \right] + k^2 \sum_{i=0}^u \frac{1}{\alpha^{2i}} \Pr(\mathcal{I}=i; u).\end{aligned}\tag{38}$$

Define the function

$$g(q, k, u) = \sum_{i=0}^u \frac{1}{\alpha^{q \cdot i}} \Pr(\mathcal{I}=i; u)\tag{39}$$

where  $\Pr(\mathcal{I}=i; u)$  is the probability distribution governing the random variable  $\mathcal{I}$ . In Section B.3 we prove Lemma B.4, which includes the following formula for  $g(2, k, u)$ :

$$g(2, k, u) = \frac{k^3 + 2k^2u + ku^2 + u(u-1)/2}{k^3}. \quad (40)$$

By the definition of  $g(q, k, u)$ , we have that Expression (38) equals:

$$\begin{aligned} & \frac{1}{1-\alpha^2} [\alpha \cdot g(2, k, u) - \alpha \cdot 1] + k^2 g(2, k, u) \\ &= \left( \frac{\alpha}{1-\alpha^2} + k^2 \right) \cdot g(2, k, u) - \frac{\alpha}{1-\alpha^2} \\ &= \frac{k(2k^2 + 2k + 1)}{2k + 1} \cdot g(2, k, u) - \frac{k(k+1)}{2k+1}. \end{aligned}$$

Finally, substituting Equation (40)'s formula for  $g(2, k, u)$  and performing elementary algebra manipulations yields the result:

$$\begin{aligned} &= \frac{2k^5 + 4k^4u + 2k^3u^2 + 3k^2u^2 + k^4 + 4k^3u + 2ku^2 + k^2u - ku + u(u-1)/2}{k^2(2k+1)} \\ &= \frac{k^2u + ku^2 + u(u-1)/2 + k^4 + 2k^3u + k^2u^2}{k^2}. \end{aligned}$$

□

### B.3 Analysis of the function $g(q, k, u)$

Recall from the proof of Lemma B.2 that  $\Pr(\mathcal{I}=i; u)$  is the probability distribution governing the final value of the Alpha Algorithm's variable  $i$ . The analysis of approximate counting in [12] includes an explanation of why the following base cases and recurrence define the distribution  $\Pr(\mathcal{I}=i; u)$ .

$$\begin{aligned} \Pr(\mathcal{I}=0; 0) &= 1 \\ \Pr(\mathcal{I}=i; 0) &= 0, \quad \forall i > 0 \\ \Pr(\mathcal{I}=0; u) &= 0, \quad \forall u > 0 \\ \Pr(\mathcal{I}=i; u) &= (1 - \alpha^i) \cdot \Pr(\mathcal{I}=i; u-1) + \alpha^{i-1} \cdot \Pr(\mathcal{I}=i-1; u-1), \quad \forall i > 0, \forall u > 0 \end{aligned}$$

Recall that  $\alpha = k/(k+1)$ , and define the function  $g(q, k, u)$  as in Equation (39):

$$g(q, k, u) = \sum_{i=0}^u \frac{1}{\alpha^{q \cdot i}} \Pr(\mathcal{I}=i; u).$$

We will now prove two lemmas that partially characterize the function  $g(q, k, u)$ . In Lemma B.3, we will prove that  $g(q, k, u)$  satisfies a certain recurrence. In Lemma B.4 we will use that recurrence to prove the correctness of explicit formulas for  $g(0, k, u)$ ,  $g(1, k, u)$ , and  $g(2, k, u)$ .

**Lemma B.3.**  $g(q, k, u)$  satisfies the following base cases and recurrence:

$$\begin{aligned} g(0, k, u) &= 1 \\ g(q, k, 0) &= 1 \\ g(q, k, u+1) &= g(q, k, u) + \left( \frac{1 - \alpha^q}{\alpha^q} \right) \cdot g(q-1, k, u) \end{aligned} \quad (41)$$

*Proof.* The base cases can be verified by inspection. The recurrence can be derived from the recurrence for  $Pr(\mathcal{I} = i; u)$  as follows:

$$\begin{aligned} g(q, k, u+1) &= \sum_{i=0}^{u+1} \frac{1}{\alpha^{q \cdot i}} \Pr(i; u+1) \\ &= \left[ \sum_{i=0}^u \frac{1}{\alpha^{q \cdot i}} (1 - \alpha^i) \Pr(i; u) \right] + 0 + 0 + \left[ \sum_{i=1}^{u+1} \frac{1}{\alpha^{q \cdot i}} (\alpha^{i-1}) \Pr(i-1; u) \right]. \end{aligned}$$

Now we will consider the two bracketed sums. For the first one:

$$\begin{aligned} &\sum_{i=0}^u \frac{1}{\alpha^{q \cdot i}} (1 - \alpha^i) \Pr(i; u) \\ &= \sum_{i=0}^u \frac{1}{\alpha^{q \cdot i}} \Pr(i; u) - \sum_{i=0}^u \frac{\alpha^i}{\alpha^{q \cdot i}} \Pr(i; u) \\ &= g(q, k, u) - \sum_{i=0}^u \frac{1}{\alpha^{(q-1) \cdot i}} \Pr(i; u) \\ &= g(q, k, u) - g(q-1, k, u). \end{aligned} \tag{42}$$

For the second one (performing the change of variables  $j = i-1$ ):

$$\begin{aligned} &\sum_{i=1}^{u+1} \frac{1}{\alpha^{q \cdot i}} (\alpha^{i-1}) \Pr(i-1; u) \\ &= \sum_{j=0}^u \frac{1}{\alpha^{q \cdot (j+1)}} (\alpha^j) \Pr(j; u) \\ &= \frac{1}{\alpha^q} \sum_{j=0}^u \frac{1}{\alpha^{(q-1)j}} \Pr(j; u) \\ &= \frac{1}{\alpha^q} g(q-1, k, u). \end{aligned} \tag{43}$$

Adding (42) and (43) yields the claimed result:

$$\begin{aligned} &g(q, k, u) - g(q-1, k, u) + \frac{1}{\alpha^q} g(q-1, k, u) \\ &= g(q, k, u) + \left( \frac{1 - \alpha^q}{\alpha^q} \right) \cdot g(q-1, k, u) \end{aligned}$$

□

**Lemma B.4.** For all integers  $u \geq 0$ :

$$\begin{aligned} g(0, k, u) &= 1 \\ g(1, k, u) &= \frac{k+u}{k} \\ g(2, k, u) &= \frac{k^3 + 2k^2u + ku^2 + u(u-1)/2}{k^3} \end{aligned}$$

*Proof.*  $g(0, k, u) = 1$  can be verified by inspection.



The result for  $q = 1$  can be proved by noting that  $\frac{k+u}{k}$  satisfies the same base cases and recurrence as  $g(1, k, u)$ . The base cases can be verified by inspection. The recurrence (41) can be verified as follows:

$$\frac{k + (u + 1)}{k} = \frac{k + u}{k} + \frac{1 - \alpha}{\alpha} \cdot 1 = \frac{k + u}{k} + \frac{1}{k}.$$

The result for  $q = 2$  can be proved by noting that  $\frac{k^3 + 2k^2u + ku^2 + u(u-1)/2}{k^3}$  satisfies the same base cases and recurrence as  $g(2, k, u)$ . The base cases can be verified by inspection. The recurrence (41) can be verified as follows:

$$\begin{aligned} & g(2, k, u) + \left( \frac{1 - \alpha^2}{\alpha^2} \right) \cdot g(1, k, u) \\ = & g(2, k, u) + \frac{1}{\alpha^2} \cdot \frac{k + u}{k} - \frac{k + u}{k} \\ = & g(2, k, u) + \frac{(k + 1)^2}{k^2} \cdot \frac{k + u}{k} - \frac{k + u}{k} \\ = & g(2, k, u) + \frac{2k^2 + k + 2uk + u}{k^3} \\ = & \frac{k^3 + 2k^2u + 2k^2 + ku^2 + 2uk + k + \frac{1}{2}u^2 + \frac{1}{2}u}{k^3} \\ = & \frac{k^3 + 2k^2(u + 1) + k(u + 1)^2 + (u + 1)u/2}{k^3} \\ = & g(2, k, u + 1). \end{aligned}$$

□

#### B.4 Analysis of the Alpha Algorithm's HIP Estimator: Proof of Theorem 4.5

*Proof.* As in Appendix B.2, let  $u = n_A - k$ , and let  $\mathcal{I}$  denote the random variable for the final value of the Alpha algorithm's variable  $i$ .

$$\begin{aligned} E(k/\alpha^{\mathcal{I}}) &= \sum_{i=0}^u \frac{k}{\alpha^i} \cdot \Pr(\mathcal{I} = i; u) \\ &= k \cdot g(1, k, u) \\ &= k \cdot \frac{k + u}{k} = (k + u) = n_A \\ \sigma^2(k/\alpha^{\mathcal{I}}) &= -E^2(k/\alpha^{\mathcal{I}}) + E((k/\alpha^{\mathcal{I}})^2) \\ &= -n_A^2 + \sum_{i=0}^u \frac{k^2}{\alpha^{2i}} \cdot \Pr(\mathcal{I} = i; u) \\ &= -n_A^2 + k^2 \cdot g(2, k, u) \\ &= \frac{u(u - 1)}{2k} \\ &= \frac{n_A^2 - 2n_Ak + k^2 - n_A + k}{2k} \\ &< \frac{n_A^2}{2k} \end{aligned}$$

□