

Methods for Finding Frequent Items in Data Streams

Graham Cormode · Marios Hadjieleftheriou

Received: date / Accepted: date

Abstract The frequent items problem is to process a stream of items and find all items occurring more than a given fraction of the time. It is one of the most heavily studied problems in data stream mining, dating back to the 1980s. Many applications rely directly or indirectly on finding the frequent items, and implementations are in use in large scale industrial systems. However, there has not been much comparison of the different methods under uniform experimental conditions. It is common to find papers touching on this topic in which important related work is mischaracterized, overlooked, or reinvented.

In this paper, we aim to present the most important algorithms for this problem in a common framework. We have created baseline implementations of the algorithms, and used these to perform a thorough experimental study of their properties. We give empirical evidence that there is considerable variation in the performance of frequent items algorithms. The best methods can be implemented to find frequent items with high accuracy using only tens of kilobytes of memory, at rates of millions of items per second on cheap modern hardware.

1 Introduction

Many data generation processes can be modeled as *data streams*. They produce huge numbers of pieces of data, each of which is simple in isolation, but which taken together lead to a complex whole. For example, the sequence of queries posed to an Internet search engine can be thought of as a stream, as can the collection of transactions across all branches of a supermarket chain. In aggregate, this data can arrive at enormous rates, easily in the realm of hundreds of gigabytes

per day or higher. While this data may be archived and indexed within a data warehouse, it is also important to process the data “as it happens”, to provide up to the minute analysis and statistics on current trends. Methods to achieve this must be quick to respond to each new piece of information, and use resources which are very small when compared to the total quantity of data.

These applications and others like them have led to the formulation of the so-called “streaming model” [38,3]. In this abstraction, algorithms take only a single pass over their input, and must accurately compute various functions while using resources (space and time per item) that are strictly sublinear in the size of the input—ideally, polynomial in the logarithm of the input size. The output must be produced at the end of the stream, or when queried on the prefix of the stream that has been observed so far. (Other variations ask for the output to be maintained continuously in the presence of updates, or on a “sliding window” of only the most recent updates). Some problems are simple in this model: for example, given a stream of transactions, finding the mean and standard deviation of the bill totals can be accomplished by retaining a few “sufficient statistics” (sum of all values, sum of squared values etc.). Others can be shown to require a large amount of information to be stored, such as determining whether a particular search query has already appeared anywhere within a large stream of queries. Determining which problems can be solved effectively within this model remains an active research area.

The frequent items problem is one of the most heavily studied questions in data streams research. The problem is popular due to its simplicity to state, and its intuitive interest and value. It is important both in itself, and as a subroutine within more advanced data stream computations. Informally, given a sequence of items, the problem is simply to find those items which occur most frequently. Typically, this is formalized as finding all items whose frequency exceeds

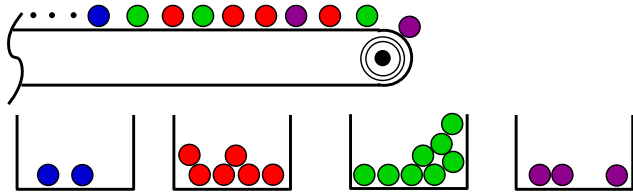


Fig. 1 A stream of items defines a frequency distribution over items. In this example, with a threshold of $\phi = 20\%$ over the 19 items grouped in bins, the problem is to find all items with frequency at least 3.8 — in this case, the green and red items (middle two bins).

a specified fraction of the total number of items. Variations arise when the items are given weights, and further when these weights can also be negative. An illustrative example is shown in Figure 1.

This abstract problem captures a wide variety of settings. The items can represent packets on the Internet, and the weights the size of the packets. Then the frequent items represent the most popular destinations, or the heaviest bandwidth users (depending on how the items are extracted from the flow identifiers). Or, the items can represent queries made to an Internet search engine, and the frequent items are now the (currently) popular terms. These are not simply hypothetical examples, but genuine cases where algorithms for this problem have been used by large corporations: AT&T [15] and Google [39] respectively. Given the size of the data (which is being generated at high speed), it is important to find algorithms which are capable of processing each new update very quickly, without blocking. It also helps if the working space of the algorithm is very small, so that the analysis can happen over many different groups in parallel, and because small structures are likely to have better cache behavior and hence further help increase the throughput.

Obtaining efficient and scalable solutions to the frequent items problem is also important since many streaming applications need to find frequent items as a ‘subroutine’ of another, more complex computation. Most directly, mining frequent *itemsets* inherently builds on finding frequent *items* as a basic building block [33]. Finding the entropy of a stream requires learning the most frequent items in order to directly compute their contribution to the entropy, and remove their contribution before approximating the entropy of the residual stream [13]. The HSS technique uses hashing to derive multiple substreams, the frequent elements of which are extracted to estimate the frequency moments of the stream [6].

Other work solves generalized versions of frequent items problems by building on algorithms for the ‘vanilla’ version of the problem. Several techniques for finding the frequent items in a sliding window model operate by keeping track of the frequent items in many sub-windows [20, 2, 32]. The ‘heavy hitters distinct’ problem, where the count of an item is the number of *distinct* pairs containing that item paired with a secondary item, is typically solved extending a fre-

quent items algorithm with distinct counting algorithms [31, 7]. Frequent items have also been applied to models of probabilistic streaming data [29], and within faster ‘skipping’ techniques [5].

Thus the problem remains an important one to understand and study in order to produce efficient streaming implementations. It remains an active area, with a steady flow of new submissions addressing the problem or variations thereof. However, sometimes prior work is overlooked or mischaracterized: algorithms first published in the eighties have been ‘rediscovered’ two decades later; existing work is sometimes claimed to be incapable of a certain guarantee, which in truth it can provide with only minor modifications; and experimental comparisons often compare against methods that are less suitable for the given problem than others that are not included (although where a subset of methods have been compared, the results are broadly in agreement with those we present here). In this paper, we try to set out clearly and concisely the main ideas in this area, as well as the common pitfalls.

Our goals are threefold:

- To provide a clear explanation of the most important algorithms for the frequent items problem, and allow comparison of their properties by using common notation and terminology. In doing so, we aim to clarify the historical development of these algorithms, and clear up some misconceptions.
- To provide baseline implementations of many of these algorithms against which future algorithms can be compared, and on top of which algorithms for different problems can be built.
- To perform a thorough experimental evaluation of the algorithms over a variety of data sets to indicate their performance in practice.

2 Definitions

Definition 1 Given a stream \mathcal{S} of n items $t_1 \dots t_n$, the frequency of an item i is $f_i = |\{j | t_j = i\}|$ (i.e. the number of indices j where the j th item is i). The exact ϕ -frequent items comprise the set $\{i | f_i > \phi n\}$.

Example. The stream $\mathcal{S} = (a, b, a, c, c, a, b, d)$ has $f_a = 3, f_b = 2, f_c = 2, f_d = 1$. For $\phi = 0.2$, the frequent items are a, b and c .

A streaming algorithm which solves this problem must use a linear amount of space, even for large values of ϕ : Given an algorithm that claims to solve this problem, we could insert a set S of N items, where every item has frequency 1. Then, we could also insert N copies of item i . If i is then reported as a frequent item (occurring more than

50% of the time) then $i \in S$, else $i \notin S$. Consequently, since correctly being able to test membership of such a set S requires $\Omega(N)$ space, $\Omega(N)$ space is also required to solve the frequent items problem. This bound holds even allowing the algorithm to use randomization. Instead, an approximate version is defined based on a tolerance for error ϵ .

Definition 2 Given a stream \mathcal{S} of n items, the ϵ -approximate frequent items problem is to return a set of items F so that for all items $i \in F$, $f_i > (\phi - \epsilon)n$, and there is no $i \notin F$ such that $f_i > \phi n$.

Since the exact ($\epsilon = 0$) frequent items problem is hard in general, we will use “frequent items” or “the frequent items problem” to refer to the ϵ -approximate frequent items problem. A related problem is to estimate the frequency of items on demand:

Definition 3 Given a stream \mathcal{S} of n items defining frequencies f_i as above, the frequency estimation problem is to process a stream so that, given any i , an \hat{f}_i is returned satisfying $\hat{f}_i \leq f_i \leq \hat{f}_i + \epsilon n$.

A solution to the frequency estimation problem allows the frequent items problem to be solved (slowly): one can estimate the frequency of every possible item i , and report those i 's whose frequency is estimated above $(\phi - \epsilon)n$. Exhaustively enumerating all items can be very time consuming (and sometimes impossible; e.g., when the items can be arbitrary strings). However, all the algorithms we study here solve both the approximate frequent items problem and the frequency estimation at the same time. Most solutions are deterministic, but we also discuss randomized solutions, which have a user-specified probability of failure.

Many other variations of the problem have been studied, and in Section 5 we discuss these and the extent to which they can be solved by extensions of the described algorithms.

3 Frequent Items Algorithms

We divide the algorithms for finding the frequent items into three classes. Counter-based algorithms track a subset of items from the inputs, and monitor counts associated with these items. For each new arrival, the algorithms decide whether to store this item or not, and if so, what counts to associate with it. A second class are derived from quantile algorithms: we show how the problem of finding (approximate) quantiles allows us to find the frequent items. Lastly, we discuss sketch algorithms, which are (randomized) linear projections of the input viewed as a vector, and solve the frequency estimation problem. They therefore do not explicitly store items from the input. In this presentation, we omit

consideration of a few algorithms based on randomly sampling items from the input, in order to keep the scope of this study bounded, and because these algorithms have attracted less interest and fewer applications.

A Note on Dictionary Issues. A common feature of several algorithms is that when given a new item, they test whether it is one of k being stored by the algorithm, and if so, increment its count. The cost of supporting this operation depends a lot on the model of computation assumed. A simple solution is to use a hash table storing the current set of items, but this means that an otherwise deterministic solution becomes randomized in its time cost, since it takes *expected* $O(1)$ operations to perform this step. Given suitable hardware, associative memory can be used to answer this in constant time; in fact, making greater use of such hardware in streaming algorithms is the subject of recent work by Bandi *et al.* [4]. But in the absence of this hardware, a dynamic dictionary data structure is needed: for example, Misra and Gries [37] discuss the use of an AVL tree. In practice, hashing is commonly used, meaning that these deterministic algorithms have randomized implementations.

3.1 Counter-based Algorithms

Majority Algorithm. The problem of frequent items dates back at least to a problem first studied by Moore in 1980. It was published as a ‘problem’ in the Journal of Algorithms in the June 1981 issue, as follows

[J.Alg 2, P208-209] Suppose we have a list of n numbers, representing the “votes” of n processors on the result of some computation. We wish to decide if there is a majority vote and what the vote is.

In addition to posing this “majority” question as a problem, Moore also invented the MAJORITY algorithm along with Boyer in 1980, described in a technical report from early 1981 [9]. To them, this was mostly of interest from the perspective of automatically proving the correctness of the solution (the details of this were published in 1991, along with a partial history [10]). In the Dec 1982 Journal of Algorithms, a solution provided by Fischer and Salzburg was published [24]. Their proposed algorithm was essentially identical to MAJORITY, although it was presented differently, and was accompanied by a proof that the number of comparisons was minimized. MAJORITY can be stated as follows: store the first item and a counter, initialized to 1. For each subsequent item, if it is the same as the currently stored item, increment the counter. If it differs, and the counter is zero, then store the new item and set the counter to 1; else, decrement the counter. After processing all items, the algorithm guarantees that if there is a majority vote, then it

Algorithm 1: FREQUENT(k)	Algorithm 2: LOSSYCOUNTING(k)	Algorithm 3: SPACESAVING(k)
<pre> 1 $n \leftarrow 0$; 2 $T \leftarrow \emptyset$; 3 foreach i do 4 $n \leftarrow n + 1$; 5 if $i \in T$ then 6 $c_i \leftarrow c_i + 1$; 7 else if $T < k - 1$ then 8 $T \leftarrow T \cup \{i\}$; 9 $c_i \leftarrow 1$; 10 else forall $j \in T$ do 11 $c_j \leftarrow c_j - 1$; 12 if $c_j = 0$ then $T \leftarrow T \setminus \{j\}$; </pre>	<pre> 1 $n \leftarrow 0$; $\Delta \leftarrow 0$; $T \leftarrow \emptyset$; 2 foreach i do 3 $n \leftarrow n + 1$; 4 if $i \in T$ then $c_i \leftarrow c_i + 1$; 5 else 6 $T \leftarrow T \cup \{i\}$; 7 $c_j \leftarrow 1 + \Delta$; 8 if $\lfloor n/k \rfloor \neq \Delta$ then 9 $\Delta \leftarrow \lfloor n/k \rfloor$; 10 forall $j \in T$ do 11 if $c_j < \Delta$ then $T \leftarrow T \setminus \{j\}$; </pre>	<pre> 1 $n \leftarrow 0$; 2 $T \leftarrow \emptyset$; 3 foreach i do 4 $n \leftarrow n + 1$; 5 if $i \in T$ then $c_i \leftarrow c_i + 1$; 6 else if $T < k$ then 7 $T \leftarrow T \cup \{i\}$; 8 $c_i \leftarrow 1$; 9 else 10 $j \leftarrow \arg \min_{j \in T} c_j$; 11 $c_i \leftarrow c_j + 1$; 12 $T \leftarrow T \cup \{i\} \setminus \{j\}$; </pre>

Fig. 2 Pseudocode for counter-based algorithms

must be the item stored by the algorithm. The correctness of this algorithm is based on a pairing argument: if every non-majority item is paired with a majority item, then there should still remain an excess of majority items. Although not posed as a streaming problem, the algorithm has a streaming flavor: it takes only one pass through the input (which can be ordered arbitrarily) to find a majority item. To verify that the stored item really is a majority, a second pass is needed to simply count the true number of occurrences of the stored item.

Frequent Algorithm. Twenty years later, two papers were published [30, 21] which include essentially the same generalization of the Majority algorithm to solve the problem of finding all items in a sequence whose frequency exceeds a $1/k$ fraction of the total count. Instead of keeping a single counter and item from the input, the FREQUENT algorithm stores $k - 1$ (item, counter) pairs. The natural generalization of the Majority algorithm is to compare each new item against the stored items T , and increment the corresponding counter if it is amongst them. Else, if there is some counter with count zero, it is allocated to the new item, and the counter set to 1. If all $k - 1$ counters are allocated to distinct items, then all are decremented by 1. A grouping argument is used to argue that any item which occurs more than n/k times must be stored by the algorithm when it terminates.

Figure 3 illustrates some of the operations on this data structure. Pseudocode to illustrate this algorithm is given in Algorithm 1, making use of set notation to represent the operations on the set of stored items T : items are added and removed from this set using set union and set subtraction respectively, and we allow ranging over the members of this set (thus implementations will have to choose appropriate data structures which allow the efficient realization of these operations). We also assume that each item j stored in T has an associated counter c_j . For items not stored in T , then c_j is defined as 0 and does not need to be explicitly stored.

It is sometimes stated that the FREQUENT algorithm does not solve the frequency estimation problem accurately, but this is erroneous. As observed by Bose *et al.*[8], executing this algorithm with $k = 1/\epsilon$ ensures that the count associated with each item on termination is at most ϵn below the true value.

The two papers published in 2002 (which cite [24]) rediscover an algorithm first published in 1982. This n/k generalization was first proposed by Misra and Gries [37]. Misra and Gries proposed a method which they refer to as “Algorithm 3”, which is equivalent to FREQUENT. In deference to this early discovery, this algorithm is also often referred to as the “Misra-Gries” algorithm in more recent work on streaming algorithms. In the same paper, the method referred to there as “Algorithm 2” correctly solves the problem but has only speculated worst case space bounds.

The time cost of the algorithm is dominated by the $O(1)$ dictionary operations per update, and the cost of decrementing counts. Misra and Gries use a balanced search tree, and argue that the decrement cost is amortized $O(1)$; Karp *et al.* propose a hash table to implement the dictionary [30]; and Demaine *et al.* show how the cost of decrementing can be made worst case $O(1)$ by representing the counts using offsets and maintaining multiple linked lists [21].

Lossy Counting. The LOSSYCOUNTING algorithm was proposed by Manku and Motwani in 2002 [33], in addition to a randomized sampling-based algorithm and techniques for extending from frequent items to frequent itemsets. The algorithm stores tuples which comprise an item, a lower bound on its count, and a ‘delta’ (Δ) value which records the difference between the upper bound and the lower bound. When processing the i th item, if it is currently stored by the algorithm in T then its lower bound is increased by one; else, a new tuple is created with the lower bound set to one, and its Δ value set to $\lfloor i/k \rfloor$. Periodically, all tuples whose upper bound is less than $\lfloor i/k \rfloor$ are deleted. These are correct upper and lower bounds on the count of each item, so at the end

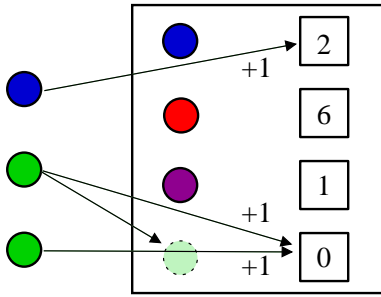


Fig. 3 Counter based data structure: the blue (top) item is already stored, so its count is incremented when it is seen. The green (middle) item takes up an unused counter, then a second occurrence increments the counter by one.

of the stream, all items whose count exceeds n/k must be stored. As with FREQUENT, setting $k = 1/\epsilon$ ensures that the error in any approximate count is at most ϵn . A careful argument demonstrates that the worst case space used by this algorithm is $O(\frac{1}{\epsilon} \log \epsilon n)$, and for certain input distributions it is $O(\frac{1}{\epsilon})$.

Storing the delta values ensures that highly frequent items which first appear early on in the stream have very accurate approximated counts. But this adds to the storage cost. A variant of this algorithm is presented by Manku in slides for the paper [34], which dispenses with explicitly storing the delta values, and instead has all items sharing an implicit value of $\Delta(i) = \lfloor i/k \rfloor$. The modified algorithm stores (item, count) pairs. For each item in the stream, if it is stored, then the count is incremented; otherwise, it is initialized with a count of 1. Every time $\Delta(i)$ increases, all counts are decremented by 1, and all items with zero count are removed from the data structure. The same proof suffices to show that the space bound is $O(\frac{1}{\epsilon} \log \epsilon n)$. This version of the algorithm is quite similar to the method presented by Misra and Gries which is labelled “Algorithm 2” in [37]; but in [34], a space bound is proven. The time cost is $O(1)$ dictionary operations, plus the periodic compress operations which require a linear scan of the stored items. This can be performed once every $O(\frac{1}{\epsilon} \log \epsilon n)$ updates, in which time the number of items stored has at most doubled, meaning that the amortized cost of compressing is $O(1)$. We give pseudocode for this version of the algorithm in Algorithm 2, where again T represents the set of currently monitored items, updated by set operations, and c_j are corresponding counts.

Space Saving. The deterministic algorithms presented thus far all have a similar flavor: a set of items and counters are kept, and various simple rules are applied when a new item arrives. The SPACESAVING algorithm of Metwally *et al.* [35] also fits this template. Here, k (item, count) pairs are stored, initialized by the first k distinct items and their exact counts. As usual, when the next item in the sequence corresponds to a monitored item, its count is incremented. But when the next item does not match a monitored item, the

(item, count) pair with the smallest count has its item value replaced with the new item, and the count incremented. So the space required is $O(k)$ (resp. $O(\frac{1}{\epsilon})$), and a short proof demonstrates that the counts of all stored items solve the frequency estimation problem with error n/k (resp. ϵn). It also shares the nice property of LOSSYCOUNTING that items which are stored by the algorithm early in the stream and not removed have very accurate estimated counts. The algorithm appears in Algorithm 3. The time cost is bounded by the dictionary operation of finding if an item is stored, and of finding and maintaining the item with minimum count. Simple heap implementations achieve this in $O(\log 1/\epsilon)$ time per update. When all updates are unitary (+1), a faster approach is to borrow ideas from the Demaine *et al.* implementation of FREQUENT, and keep the items in groups with equal counts. By tracking a pointer to the group with smallest count, the find minimum operation takes constant time, while incrementing counts takes $O(1)$ pointer operations (the “Stream-Summary” data structure in [35]).

3.2 Quantile Algorithms

The problem of finding the ϕ -quantiles of a sequence of items drawn from a totally ordered domain is to find an item i such that it is the smallest item which dominates ϕn items from the input. We define the *rank* of item i as $\text{rank}(i) = \sum_{j < i} f_j$; i.e., this is the total frequency of all items dominated by i . So the ϕ quantile is the i which satisfies $\text{rank}(i) \leq \phi n$ and $\text{rank}(i+1) > \phi n$. The approximate version allows ϵn uncertainty in the ranks, i.e., to find an i such that $\text{rank}(i) \leq (\phi + \epsilon)n$ and $\text{rank}(i+1) > (\phi - \epsilon)n$.

This problem is more general than frequent items, since a solution to the approximate quantiles problem allows frequent items to be found, by the following observation: suppose i is a frequent item with $f_i > 2\epsilon n$. Then i *must* be reported as the approximate ϕ quantile for all values of ϕ in the range $\text{rank}(i) + \epsilon$ to $\text{rank}(i+1) - \epsilon$, and by our assumption on f_i , this range is non-empty. Similarly, if the quantile algorithm also produces an estimate of $\text{rank}(i)$ with error at most ϵn , this can be used to solve the frequency estimation problem, since $f(i) = \text{rank}(i+1) - \text{rank}(i)$.

GK Algorithm. The approximate quantiles algorithm of Greenwald and Khanna [26], usually referred to simply as the GK algorithm is somewhat similar to LOSSYCOUNTING, in that it stores tuples containing an item from the input, a frequency count g , and a Δ value. Here though, the tuples are kept sorted under the total order of the domain of items. The g value encodes the difference between the lowest possible rank of the stored item and the previous stored item; the Δ value encodes the difference between the greatest possible rank of the item and lowest possible rank. An estimated rank of any item (whether it is stored by the algorithm or not) can

be computed from this information. Every new arrival is inserted as a new tuple in the sorted order with a g value of 1 and a Δ value of $\lfloor \epsilon n \rfloor$. This ensures that the requirements on g and Δ are met. Periodically, a “compress” operation removes some tuples: for two adjacent tuples i and $i + 1$, if $g_i + g_{i+1} + \Delta_{i+1} \leq \epsilon n$, then the i th tuple is removed, and we set $g_{i+1} \leftarrow g_{i+1} + g_i$. It can be shown that this allows the rank of any item to be estimated with error at most ϵn , and that (under a slightly formalized version of the algorithm) the space required is bounded by $O(\frac{1}{\epsilon} \log \epsilon n)$. The time cost requires inserting new tuples into a list in sorted order, and periodically scanning this list to merge some adjacent tuples. This can be supported in (amortized) time logarithmic in the size of the data structure.

QDigest. The QDIGEST algorithm was proposed by Suri *et al.* [41] in the context of monitoring distributed data. However, it naturally applies to a streaming setting. It assumes that the ordered domain can be represented as the set of integers $\{1 \dots U\}$. Each tuple stored by QDIGEST consists of a dyadic range and a count. A dyadic range is a range whose length is a power of two, and which ends at a multiple of its own length, i.e., can be written as $\{j2^\ell - 1 \dots (j + 1)2^\ell\}$. Each new item i is inserted as a trivial dyadic range $\{i\}$ with count 1 (or the count of that range is incremented if it is already present in the data structure). Observe that every non-trivial dyadic range can be partitioned into two dyadic ranges of half the length; call the set comprised of a range and its two half-length subranges a triad. The algorithm enforces the invariants that each non-trivial range has an associated count at most $\frac{\epsilon n}{\log U}$; and that the sum of counts associated with every triad is at least $\frac{\epsilon n}{\log U}$. If the second of these does not hold, then a compressing operation removes the counts associated with the two subranges and adds them on to the parent range. It is then straightforward to show that when these invariants hold, the total number of ranges with non-zero counts is at most $O(\frac{\log U}{\epsilon})$. Moreover, the true frequency of an item is at most the sum of counts of the $\log U$ ranges which contain that item, and since their count is bounded, the count of the trivial range corresponding to an item is at most an ϵn underestimate. Hence, the data structure directly solves the frequency estimation problem, and further, a walk over the induced tree structure in time linear in the data structure size extracts all frequent items.

In implementing this algorithm, the main challenge is to implement the compress operation so that its cost is minimized and it correctly restores the required invariants on execution: it is not clear that the version of compress presented in the original q-digest paper [41] restores the invariant. Subsequently, alternate implementations of the q-digest, with additional properties (such as ensuring that if a range is present in the summary, then so is its parent range) have been proposed [28, 16]. The time cost of the best of these is (amortized) $O(\log \log U)$.

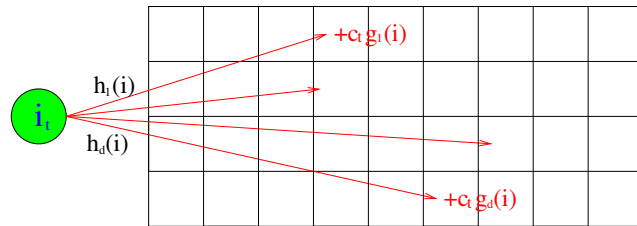


Fig. 4 Sketch data structure: each new item is mapped to a set of counters, which are incremented.

3.3 Sketches

Here, we use the term ‘sketch’ to denote a data structure which can be thought of as a linear projection of the input. That is, if we imagine the stream as implicitly defining a vector whose i -th entry is f_i , the sketch is the product of this vector with a matrix. For the algorithm to use small space, this matrix will be implicitly defined by a small number of bits. The algorithms use hash functions to define the linear projection. There is sometimes confusion on this issue, but it is straightforward to interpret the algorithms below which are defined in terms of using hash functions to map items to array entries as also defining a (sparse) matrix. Hence, it is meaningful to use both hashing and linear projection terminology to talk about sketches, and there is no need to draw a distinction between the two perspectives.

The sketch algorithms solve the frequency estimation problem, and so need additional information to solve the frequent items problem. We outline two sketching approaches below, followed by methods which augment the stored sketch to find frequent items quickly.

CountSketch. The first sketch in the sense that we use the term was the AMS or Tug-of-war sketch due to Alon *et al.* [1]. This was used to estimate the second frequency moment, $F_2 = \sum_i f_i^2$. It was subsequently observed that the same data structure could be used to estimate the inner-product of two frequency distributions, i.e., $\sum_i f_i f'_i$, for two distributions given (in a stream) by f_i and f'_i . But this means that if f_i is defined by a stream, at query time we could find the product with $f'_i = 1$ and $f'_j = 0$ for all $j \neq i$. Then, the true answer to the inner product should be exactly f_i . The error guaranteed by the sketch turns out to be $\epsilon F_2^{1/2} \leq \epsilon n$ with probability at least $1 - \delta$ for a sketch of size $O(\frac{1}{\epsilon^2} \log 1/\delta)$. The ostensibly dissimilar technique of “Random Subset Sums” [25] (on close inspection) turns out to be isomorphic to this instance of the algorithm.

Maintaining this data structure is slow, since it requires updating the whole sketch for every new item in the stream. The COUNTSKETCH algorithm of Charikar *et al.* [14] dramatically improves the speed by showing that the same underlying technique works if each update only affects a small subset of the sketch, instead of the entire summary. The

Algorithm 4: COUNTSKETCH(w, d)

```

1  $C[1, 1] \dots C[d, w] \leftarrow 0;$ 
2 for  $j \leftarrow 1$  to  $d$  do
3    $\lfloor$  Initialize  $g_j, h_j;$ 
4 foreach  $i$  do
5    $n \leftarrow n + 1;$ 
6   for  $j \leftarrow 1$  to  $d$  do
7      $C[j, h_j(i)] \leftarrow C[j, h_j(i), j] + g_j(i);$ 

```

Algorithm 5: COUNTMIN(w, d)

```

1  $C[1, 1] \dots C[d, w] \leftarrow 0;$ 
2 for  $j \leftarrow 1$  to  $d$  do
3    $\lfloor$  Initialize  $g_j;$ 
4 foreach  $i$  do
5    $n \leftarrow n + 1;$ 
6   for  $j \leftarrow 1$  to  $d$  do
7      $C[j, h_j(i)] \leftarrow C[j, h_j(i)] + 1;$ 

```

Fig. 5 Pseudocode for sketching algorithms

sketch consists of a $d \times w$ array C of counters, and two hash functions for each of the d rows, h_j which maps input items onto $[w]$, and g which maps input items onto $\{-1, +1\}$. Each input item i causes $g_j(i)$ to be added on to entry $C[j, h_j(i)]$ in row j , for $1 \leq j \leq d$. The estimate \hat{f}_i is $\text{median}_{1 \leq j \leq d} g_j(i) C[j, h_j(i)]$. The estimate derived for each value of j can be shown to be correct in expectation and has variance depending on F_2/w . Figure 4 shows a schematic of the data structure under the update procedure: the new item i gets mapped to a different location in each row, where $g_j(i)$ is added on to the current counter value in that location. We illustrate the core of the update algorithm in Algorithm 4. Using d rows drives down the probability of giving a bad estimate, so setting $d = \log \frac{4}{\delta}$ and $w = O(\frac{1}{\epsilon^2})$ ensures that \hat{f}_i has error at most $\epsilon F_2^{1/2} \leq \epsilon n$ with probability at least $1 - \delta$. Giving this guarantee requires that each g and h is picked independently from a family of pairwise independent hash functions. Efficient implementations of such hash functions are described by Thorup [42]. The total space used is $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$, and the time per update is $O(\log \frac{1}{\delta})$ worst-case.

CountMin Sketch. The COUNTMIN sketch algorithm of Cormode and Muthukrishnan [19] can be described in similar terms to COUNTSKETCH. As before, an array of $d \times w$ counters is maintained, and pairwise independent hash functions h_j map items onto $[w]$ for each row. Each update is mapped onto d entries in the array, each of which is incremented. Now $\hat{f}_i = \min_{1 \leq j \leq d} C[j, h_j(i)]$. The Markov inequality is used to show that the estimate for each j overestimates by less than n/w , and repeating d times reduces the probability of error exponentially. So setting $d = \log \frac{1}{\delta}$ and $w = O(\frac{1}{\epsilon})$ ensures that \hat{f}_i has error at most ϵn with probability at least $1 - \delta$. Consequently, the space is $O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ and the time per update is $O(\log \frac{1}{\delta})$. The data structure and update procedure is consequently much like that illustrated for the Count-Sketch in Figure 4, with $g_j(i)$ always equal to 1. The update algorithm is shown in Algorithm 5.

Finding Frequent Items using a Hierarchy. Sketches allow us to model the removal of items (to denote the conclusion of a packet flow; or the return of a previously bought item, say) as an update with negative weight. Two variations follow: the “strict” version, where the input is guar-

anteed at all times to induce non-negative frequencies, and the “general” case, where the total weight of an item is allowed to be negative. Only sketch based algorithms have been shown to adapt to these settings. In the strict case, an approach based on divide-and-conquer will work: additional sketches are used to determine which (dyadic) ranges of items are frequent [19]. If a range is frequent, then it can be split into two equal subranges, and the frequency of each subrange estimated from an appropriate sketch, until a single item is returned. More generally, rather than splitting the range into 2 subranges, we can split into b . This trades off update time against query time: if all items $i \in \{1 \dots U\}$, then $\lceil \log_b U \rceil$ sketches suffice, but each potential range is split into $b > 1$ subranges when answering queries. Thus, updates take $O(\log_b U \log \frac{1}{\delta})$ hashing operations, and $O(1)$ counter updates for each hash. Typically, moderate constant values of b are used (between 2 and 256, say); choosing b to be a power of two allows fast bit-shifts to be used in query and update operations instead of slower divide and modulo operations. This results in COUNTMIN sketch Hierarchical and COUNTSKETCH Hierarchical algorithms.

Finding Frequent Items using Group Testing. In the general case, even this fails, and new techniques are needed [18, 40]. The idea of “group testing” in this context [18] randomly divides the input into buckets so that we expect at most one frequent item in each group. Within each bucket, the items are divided into groups so that the “weight” of each group indicates the identity of the frequent item. This can be seen as an extension of the Count-Min sketch, since the structure resembles the buckets of the sketch, with additional information on subgroups of each bucket (based on the binary representation of items falling in the bucket); further, the analysis and properties are quite close to those of a Hierarchical Count-Min sketch. For each bucket, we keep additional counts for the total frequency of all items whose binary representation has the i th bit set to 1. This increases the space to $O(\frac{1}{\epsilon} \log U \log \delta)$ when the binary representation takes $\log U$ bits. Each update requires $O(\log \frac{1}{\delta})$ hashes as before, and updating $O(\log U)$ counters per hash.

4 Experiments

4.1 Setup

We ran several algorithms under a common implementation framework to test as accurately as possible their relative performance. All algorithms were implemented using C++, and used common subroutines for similar tasks (e.g., hash tables) to increase comparability. We ran experiments on a 4 Dual Core Intel(R) Xeon(R) 2.66 GHz with 16 GB of RAM running Windows 2003 Server. The code was compiled using Microsoft's Visual C++ 2005 compiler and g++ 3.4.4 on cygwin. We did not observe significant differences between the two compilers. We report here results obtained using Visual C++ 2005. The code extended and enhanced the MassDal implementations <http://www.cs.rutgers.edu/~muthu/massdal-code-index.html>; the new versions can be downloaded from <http://www.research.att.com/~marioh/frequent-items>.

For every algorithm we tested a number of implementations, using different data structures to implement the basic set operations. For some algorithms the most robust implementation was obvious. For other algorithms we present here results of competing solutions. We compare counter based algorithms, quantile estimation algorithms, and sketch algorithms. For counter based algorithms we examine: FREQUENT using the Demaine *et al* implementation technique of linked lists (F), LOSSYCOUNTING keeping separate delta values for each item (LCD), LOSSYCOUNTING without deltas (LC), SPACESAVING using a heap (SSH), and SPACESAVING using linked lists (SSL). For quantile algorithms we examine: GK (GK) and QDIGEST (QD). Finally, we examine the following sketches: hierarchical COUNTSKETCH (CS), hierarchical COUNTMIN sketch (CMH), and the Combinatorial Group Testing variant of COUNTMIN (CGT). We separate these comparisons into the three categories of algorithms, since each group has different characteristics. Counter based algorithms solve only the frequent elements problem. Quantile algorithms are also good for estimating quantiles, and hence more powerful. Sketches work under both insertions and deletions and are the only alternative in applications that need to support deletions. The added functionality of quantile and sketch algorithms comes at a cost; usually, either at the expense of reduced update throughput, or increased memory consumption. Previous work has not distinguished these classes, leading to the observation that sketch algorithms require more space than counter-based algorithms, although the classes really apply to different scenarios.

We ran experiments using real network traffic and generated data. The network data set was drawn from 24 hours of traffic from a router in a major network backbone. We ran experiments using 10 million packets of HTTP traffic,

and 10 million packets of all UDP traffic. We generated data from a skewed distribution (Zipf), varying the skew from 0.8 to 2 (in order to obtain meaningful distributions that produce at least one heavy hitter per run). Finally, we also varied the frequency threshold ϕ , from 0.0001 to 0.01. In our experiments, we set the error guarantee $\epsilon = \phi$, since our results showed that this was sufficient to give high accuracy in practice.

We compare the efficiency of the algorithms with respect to:

- Update throughput, measured in number of updates per millisecond.
- Space consumed, measured in bytes.
- Recall, measured in the total number of true heavy hitters reported over the number of true heavy hitters given by an exact algorithm.
- Precision, measured in total number of true heavy hitters reported over the total number of answers reported. Precision quantifies the number of false positives reported.
- Average relative error of the reported frequencies. We measure separately the average relative error of the frequencies of the true heavy hitters, and the average relative error of the frequencies of the false positive answers.

For all of the above, we perform 20 runs per experiment (by dividing the input data into 20 chunks and querying the algorithms once at the end of each run). Furthermore, we ran each algorithm independently from the others to take advantage of possible caching effects. We report averages on all graphs, along with the 5th and 95th percentiles as error bars.

4.2 Counter based algorithms

In this section we compare FREQUENT (F), LOSSYCOUNTING with delta values (LCD), LOSSYCOUNTING without deltas (LC), SPACESAVING using a heap (SSH), and SPACESAVING using linked lists (SSL). First we present results for the Zipf generated data. The default skew parameter, unless otherwise noted, is $z = 1.0$, and the default frequency threshold is $\phi = 0.001$. Then, we show trends for the network traffic data.

Space and Time costs. Figures 6(a) and 6(b) show the update throughput of the algorithms as a function of data skew (z) and increasing frequency threshold (ϕ) respectively. We can see that update throughput increases significantly for highly skewed data. This is expected, since high skew translates to a very small number of truly frequent items, simplifying the problem. SSL is very fast, but SSH (the same algorithm, but implemented with a heap) is appreciably slower. This shows how data structure choices can affect the performance. In fact, one might expect the heap based solution to

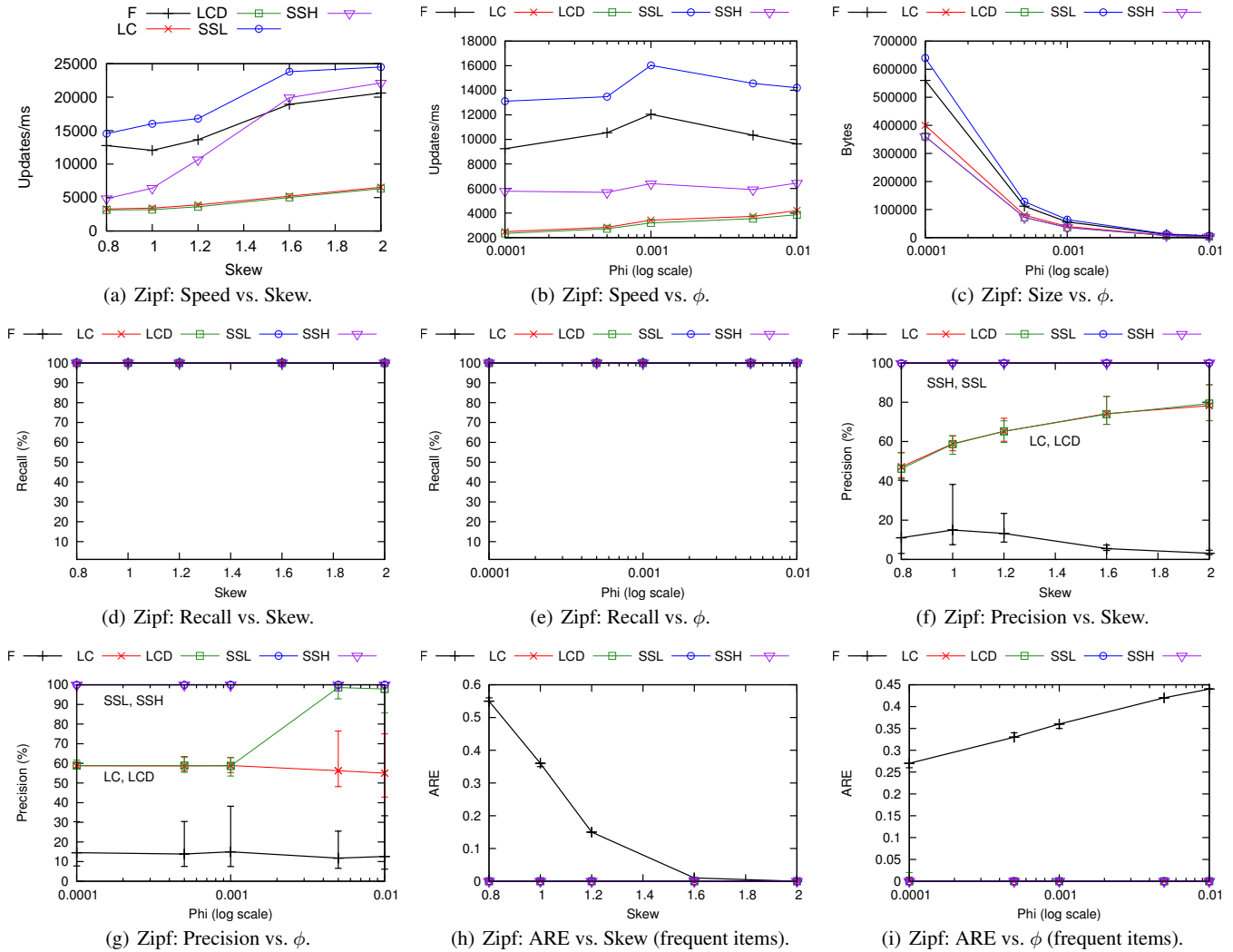


Fig. 6 Performance of counter-based algorithms on synthetic data

be even slower. However, since SSH has to restore the min heap property after each update, in practice, updates only slightly change the frequency of an item, so this only takes a small number of operations. It is also clear that the range of frequency thresholds (ϕ) considered did not affect update throughput (notice the log scale on the horizontal axis). As we see in the subsequent plots, the summary structures fit within a modern second level cache, so there is no obvious effect due to crossing memory boundaries here. We would expect smaller thresholds to result in a larger number of frequent elements, hence larger data structures and decreased throughput. In our datasets, the total number of truly frequent elements for $z = 1.0$ varies from 713 to 6 as ϕ increases; a small enough difference not to affect running times.

Figure 6(c) plots the space consumed by each structure. In our implementations of LOSSYCOUNTING, the maximum number of counters was fixed as a function of ϕ , to avoid

memory allocation during stream processing. For the other algorithms, the space used is directly determined by ϕ . So the space consumed is not affected by skewness for fixed ϕ (hence we omit the plot). Varying ϕ has a direct effect. Smaller ϕ 's imply a significantly larger number of candidates exceeding the frequency threshold that need to be maintained. It should be noted here that, for our datasets, a naive solution that maintains one counter per input item would consume many megabytes (and this grows linearly with the input size). This is at least 12 times larger than SSH for $\phi = 0.0001$ (which is the most robust algorithm in terms of space), and over a thousand times larger than all algorithms for $\phi = 0.01$. Clearly, the space benefit of these algorithms even for small frequency thresholds is substantial in practice.

Precision and Recall. Figures 6(d) and 6(e) plot recall, computed as the total number of true frequent items returned over the exact number of frequent items. The deterministic

Algo/Skew	0.8	1.0	1.2	1.6	2.0
F	0.84	0.84	0.80	0.63	0.40
LC	11.29	6.93	2.13	2.16	1.52
LCD	15.26	10.52	7.56	1.86	1.29
SSL	0	0	0	0	0
SSH	0	0	0	0	0

(a) Zipf: ARE vs Skew for false positives.

Algo/ ϕ	0.0001	0.005	0.001	0.05	0.01
F	0.86	0.86	0.84	0.81	0.77
LC	3.97	6.37	6.93	4.78	3.40
LCD	0	0	10.52	5.30	3.84
SSL	0	0	0	0	0
SSH	0	0	0	0	0

(b) Zipf: ARE vs ϕ for false positives.

Fig. 7 Average Relative Error (ARE) on Zipfian Data

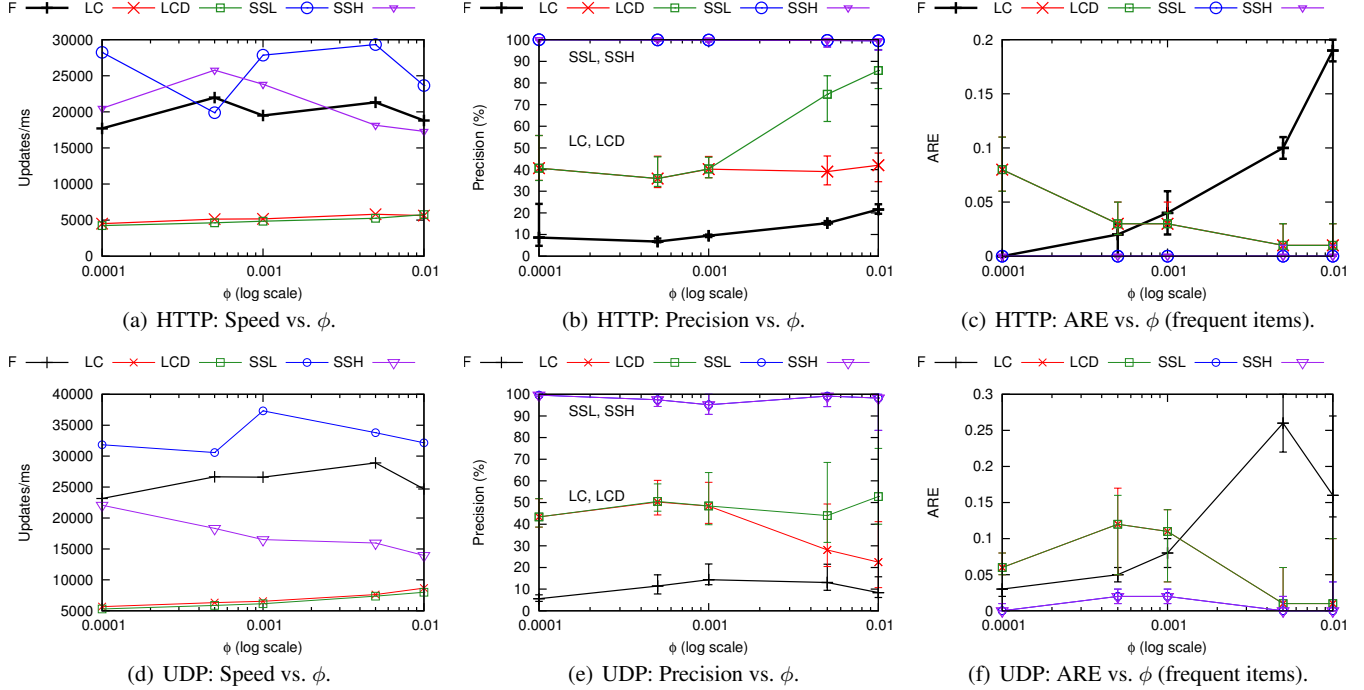


Fig. 8 Performance of counter-based algorithms on real network data

algorithms guarantee to return all ϕ frequent items, and possibly some false positives, so we expect 100% recall, which is observed in the plots. Figures 6(f) and 6(g) plot precision. We also show the 5th and 95th percentiles in the graphs as error bars. Precision is the total number of true answers returned over the total number of answers. Precision is an indication of the number of false positives returned. Higher precision means smaller number of false positive answers. There is a clear distinction between different algorithms in this case. When using $\epsilon = \phi$, F results in a very large number of false positive answers, while LC and LCD result in approximately 50% false positives for small skew parameters, but their precision improves as skewness increases. This is expected since frequent items are easier to identify for highly skewed data where the number of potentially frequent candidates is small, and there are fewer “almost frequent” items. Decreasing ϵ relative to ϕ would improve this at the cost of increasing the space used. However, SSL and SSH yield 100% accuracy in all cases (i.e., no false positives), with about the same or better space usage. Note that

these implement the same algorithm and so have the same output, only differing in the underlying implementation of certain data structures. Finally, notice that by keeping additional per-item information, LCD distinguishes between truly frequent and potentially frequent items marginally better than LC.

Relative Error. Figures 6(h) and 6(i) plot the average relative error in the frequency estimation of the truly frequent items. The graph also plots the 5th and 95th percentiles as error bars. All algorithms except F have zero estimation error with zero variance. Clearly sophisticated counter based algorithms are able to track the exact frequency of the truly frequent items exactly, which is expected. F yields very large frequency estimation errors for low skew, but the error drops as the skew increases. The variance is very small in all cases. On the other hand, estimation error for F increases as ϕ increases.

Figures 7(a) and 7(b) show the average relative error in the frequency estimation of false positive answers. SSL and SSH do not report any false positives for average and high

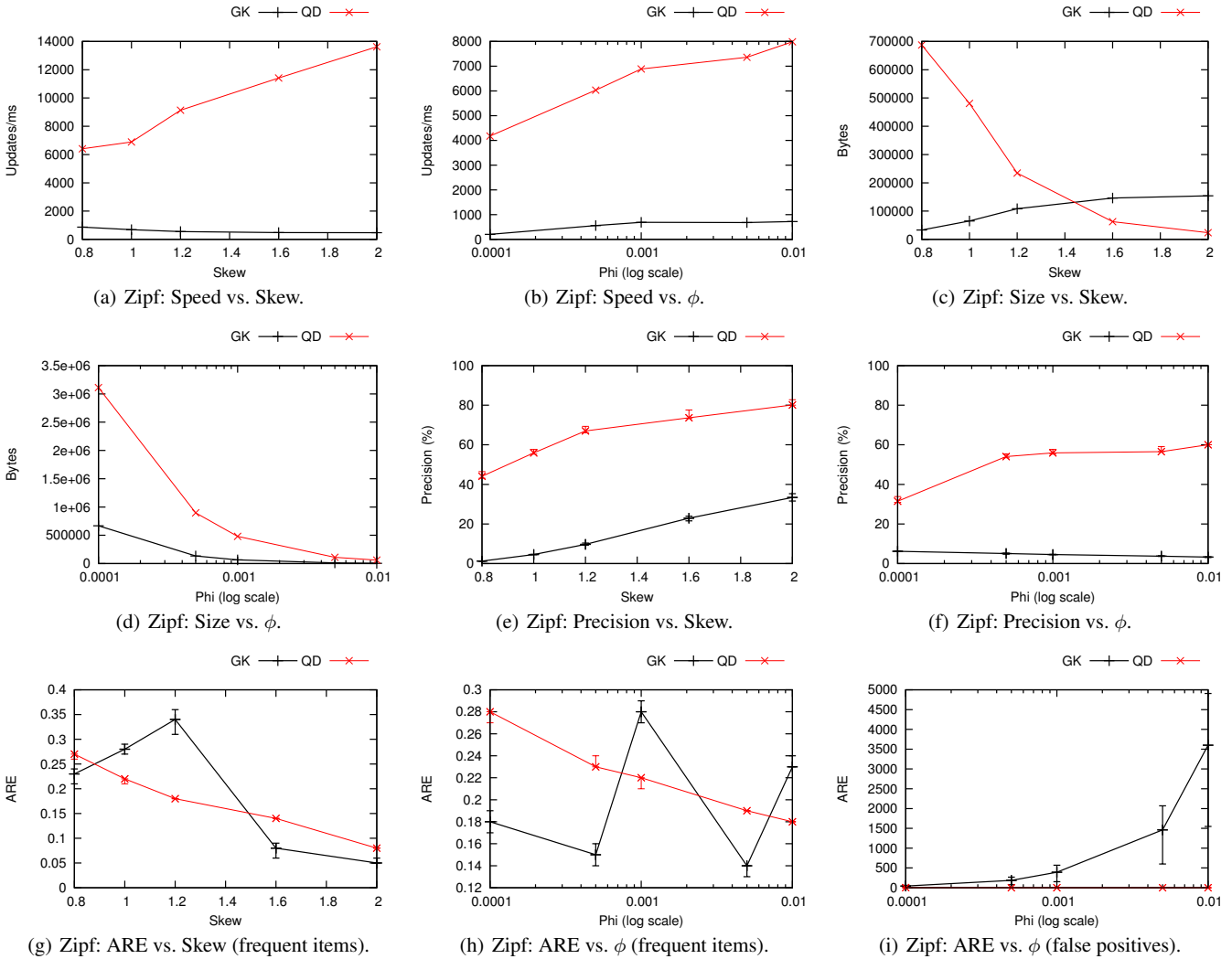


Fig. 9 Performance of quantile algorithms on synthetic data

skew, hence the error is zero. For the rest of the algorithms it is clear that the estimated frequencies of non-frequent items can be far from their true values. F always returns an underestimate of the true count of any item, hence its errors are less than 1; LC and LCD always return overestimates based on a Δ value, and so yield inflated estimates of the frequencies of infrequent items.

Network Data. Finally we ran the same experiments on real network traffic. Figures 8(a) to 8(c) show results for HTTP traffic, while Figures 8(d) to 8(f) for UDP traffic. In both cases, we track the most frequent destination IP addresses. We plot everything as a function of ϕ . The trends observed are similar to the ones for generated data, hence we omit a more detailed analysis for brevity.

Conclusion. Overall, the SPACESAVING algorithm appears conclusively better than other counter-based algorithms, across a wide range of data types and parameters. Of the two implementations compared, SSH exhibits very good performance

in practice. It yields very good estimates, with 100% recall and precision, consumes very small space and is fairly fast to update (faster than LC and LCD). Alternatively, SSL is the fastest algorithm with all the good characteristics of SSH, but consumes twice as much space on average. If space is not a critical issue SSL is the implementation of choice.

4.3 Quantile algorithms

Quantile structures are more expensive to update and store compared to counter based algorithms, but they solve a more general problem. When a quantile estimation algorithm needs to be maintained, it can be used to solve the frequent items problem as well. In this section we compare the GK and QDigest algorithms. We run the same set of experiments, using a default value of $z = 1.0$ and $\phi = 0.001$.

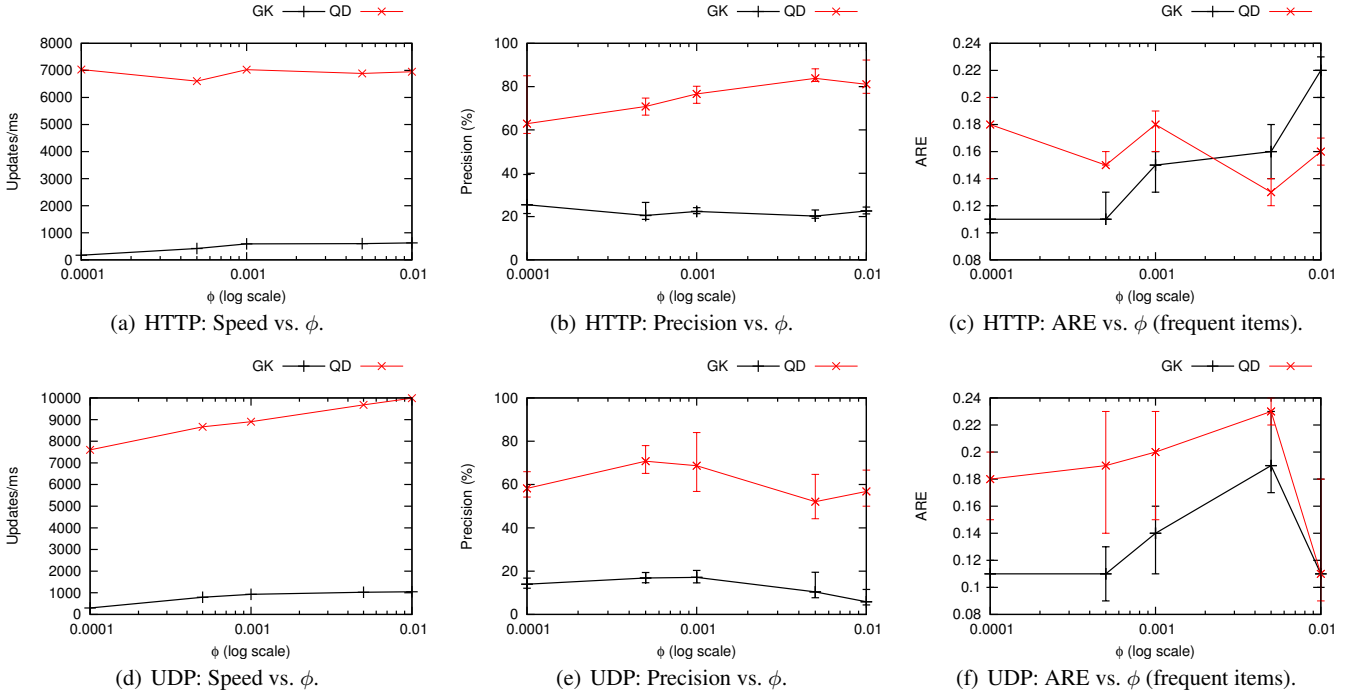


Fig. 10 Performance of quantile algorithms on real network data

Space and Time Costs. Figures 9(a) and 9(b) show the update throughput of the algorithms. GK is not affected by data skewness, while QD becomes faster as the data becomes more skewed. Increasing frequency thresholds (ϕ) has a positive effect on update performance, especially for QD probably due to the reduced structure size (from 4MB to less than 100KB as ϕ varies). Figures 9(c) and 9(d) plot the space consumed. Notice that data skewness affects the structure size of the quantile algorithms (in contrast with counting based algorithms). The QD algorithm is able to compress into a smaller data structure when the data are more skewed, since more of the total “weight” of the input is stored in a small number of leaves in the tree. On the other hand the GK algorithm is negatively affected as skew increases. Notice that the size of the quantile structures is up to 7 times larger than the most space inefficient counter based algorithm. Analytically, this cost is a logarithmic factor ($\log U$ for QD, $\log(\epsilon n)$ for GK), which seems to be the root cause of the higher cost. Indeed, for small enough values of ϕ , QD begins to approach the size of the naive solution.

Precision, Recall, and Error. We omit the figures for recall, since both algorithms have 100% recall in all cases. Figures 9(e) and 9(f) plot precision. The precision of GK is very low. The algorithm reports a large number of false positives. The precision of QD improves as the skew increases and as ϕ increases, but remains below 80% in all cases. This is not a surprising result, since these algorithms are not tailored for frequent item identification, and are being run with ϵ being

large relative to ϕ ; reducing ϵ would improve precision, but cost yet more in space.

Figures 9(g) and 9(h) plot the average relative error in the frequency estimation of the truly frequent items. These quantile estimation algorithms are not very accurate for frequency estimation, for average data skewness. Finally, Figure 9(i) plots the average relative error in the frequency estimation of false positives. Here, GK vastly overestimates the frequency of some rare items, and QD also inflates frequencies somewhat. Straightforwardly, we do not expect the algorithms to give useful results without increasing the already high space costs, due to the very large number of false positive answers and the fact that the algorithms are not robust in estimating frequencies in general.

Network data. The performance of quantile algorithms on network data is shown in Figure 10. The trends are the same as for synthetic data, both for HTTP and UDP traffic, with the only observable difference being slightly larger precision and smaller relative errors.

Conclusion. Overall, the quantile algorithms cannot compete against the counter based algorithms for identifying heavy hitters. They have larger structures, are slower to update, and yet still do not estimate frequencies accurately.

4.4 Sketch algorithms

Finally, we evaluate sketching algorithms. The advantage of sketches is that they support deletions, and hence are the

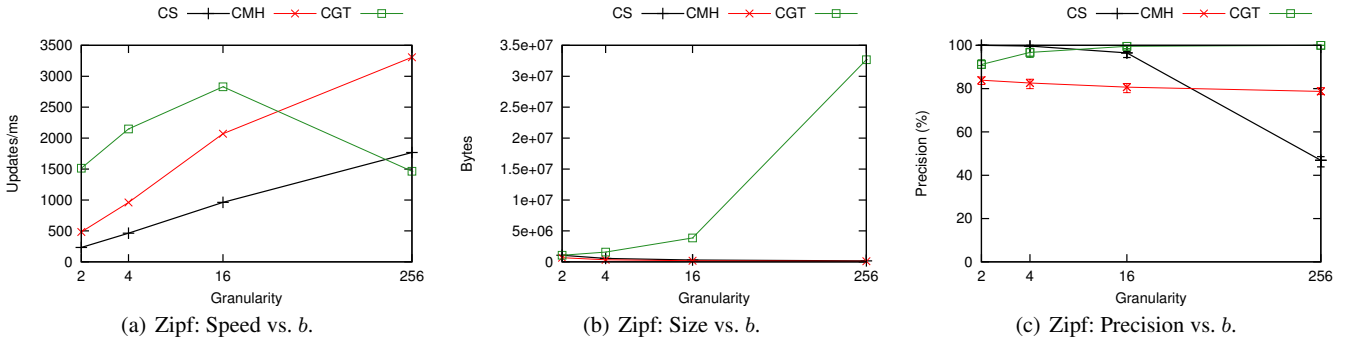


Fig. 11 Performance of sketch algorithms on synthetic data, for various granularities

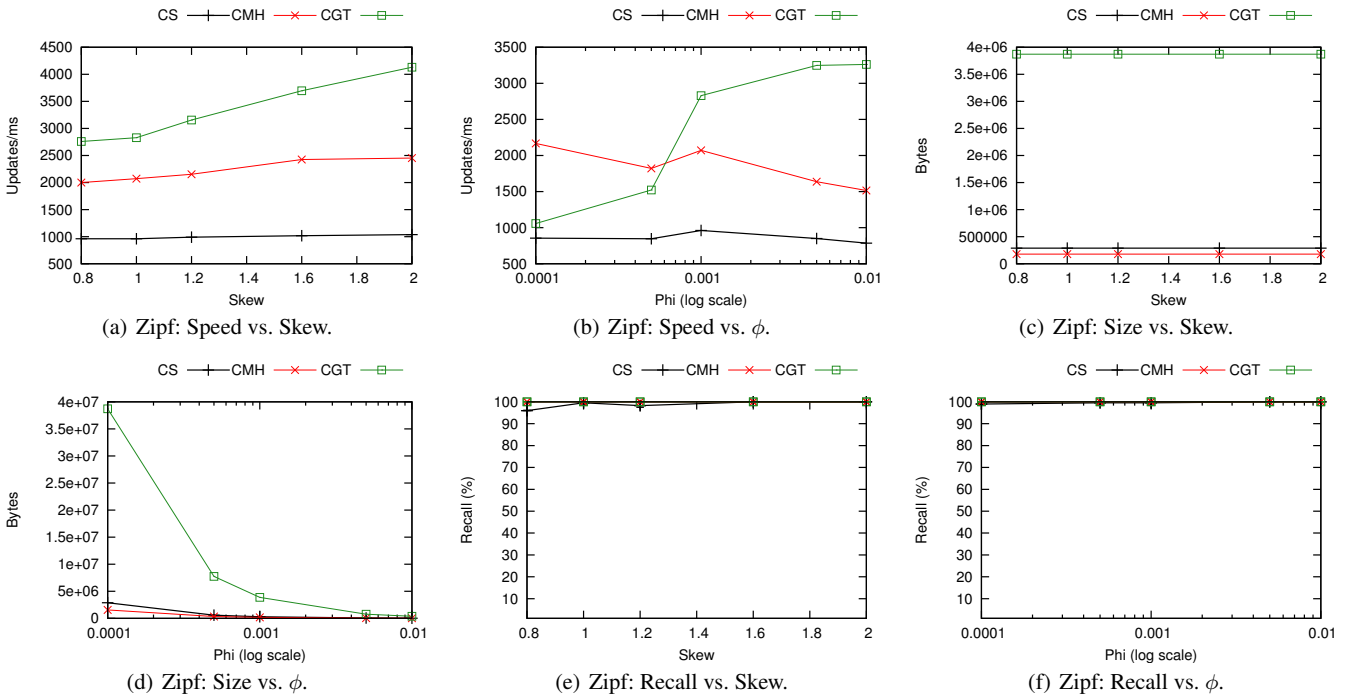


Fig. 12 Performance of sketch algorithms on synthetic data (Speed, size and recall)

only alternative in fully dynamic environments (that is, cases where there is a large number of deletions and insertions). This comes at the cost of increased space consumption and slower update performance. We run the same set of experiments, using a default value of $z = 1.0$ and $\phi = 0.001$. Recall that the algorithms all take a parameter b which determines the “fan-out” of a hierarchy used to search for the frequent items. This parameter b is usually chosen to be a power of 2 for efficiency. We compared various values of b , as shown in Figure 11. Subsequently, we use a hierarchy with $b = 16$ for all algorithms, since this seems to offer the the best tradeoff between speed, size and precision. The sketch depth is set to $d = 4$ throughout, and the width to $w = 2/\phi$, based on the analysis of the COUNTMIN sketch. This keeps the space of CS and CMH relatively close, and CGT constant factors larger.

Space and Time Cost. Figures 12(a) and 12(b) show the update throughput of the algorithms. Update throughput is not affected by data skewness, and marginally affected by variations in ϕ , except for the CGT algorithm. CS has the slowest update rate among all algorithms, due to the larger number of hashing operations needed. The fastest sketch algorithm is from 5 up to 10 times slower than the fastest counter based algorithm. Figures 12(c) and 12(d) plot the space consumed. The size of the sketches is fairly large compared to counter based algorithms. CMH is the most space efficient sketch and still consumes space 3 times as large as the least space efficient counter based algorithm.

Precision, Recall and Error. Figures 12(e) and 12(f) plot recall. We observe that for the sketches the recall is not always 100%. The error of CMH is one sided, and as a consequence, it still guarantees 100% recall; CGT does not have

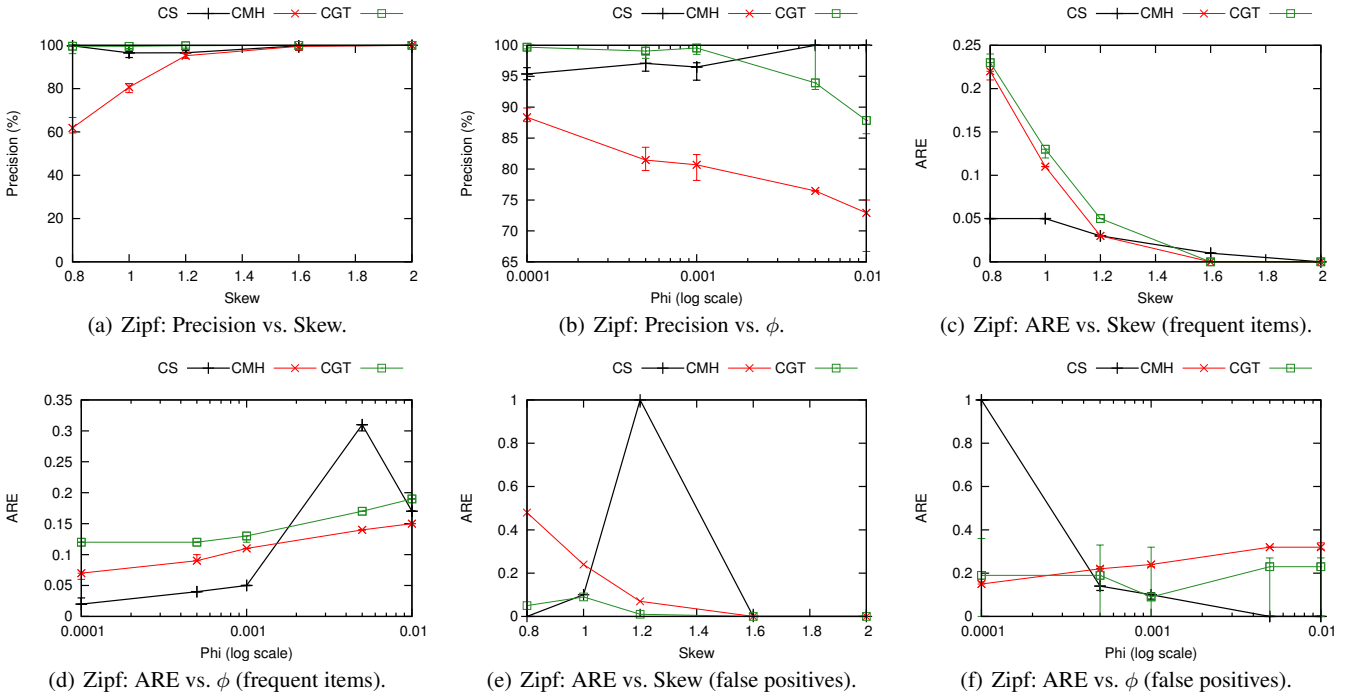


Fig. 13 Performance of sketch algorithms on synthetic data (precision and average relative error)

as strong a guarantee, but also achieved 100% recall in all our experiments. CS has a higher probability of failing to recover some frequent items, but still achieved close to 100% in all cases (94% in the worst case). Since these are probabilistic algorithms, in practice there is a small probability of failure (equal to 1% for our experiments). Nevertheless, overall recall is very close to 100% in all cases.

Figures 13(a) and 13(b) plot precision. CMH has low precision for average skewness, but improves as data skew increases. CMH can accurately estimate the frequencies of only the most frequent items. As the distribution approaches a uniform distribution the error in estimated frequencies increases, and the number of false positives increases. The other algorithms exhibit precision higher than 85% in all cases.

Figures 13(c) and 13(d) plot the average relative error in the frequency estimation of the truly frequent items. For sufficiently skewed distributions all algorithms can estimate item frequencies very accurately. Results here are very comparable, since they essentially correspond to a single instance of a COUNTSKETCH or COUNTMIN sketch, both with the same amount of space for each point plotted. Hence CMH and CGT are quite similar, corresponding to a single COUNTMIN sketch (with different random choice of hash functions, leading to variations). We do not expect ϕ to affect estimation accuracy significantly, but as ϕ increases the total number of answers decreases which could make a difference. CS exhibits erratic behaviour here. This could be attributed to a

random failure of the sketch, since this is a probabilistic algorithm¹.

Finally, Figures 13(e) and 13(f) plot the average relative error in the frequency estimation of false positives. The errors increase substantially in this case for CMH and low skew data. CS seems to be able to estimate frequencies very accurately, but exhibits outliers once again.

Network data. Figure 14 shows results for network data. Once again, recall for all algorithms was almost 100% and is not shown in the figure. The general trends here are almost the same as for synthetic datasets, the only exception being the CS algorithm having lower precision for HTTP data. Nevertheless, precision is once again higher than 75% in all cases. Notice also that the average relative error of frequencies is significantly smaller for UDP data. It is clear that these algorithms are affected heavily by the application and data at hand.

Conclusion. There is no clear winner among these algorithms. CMH has small size and high update throughput, but is only accurate for highly skewed distributions. CGT consumes a lot of space but it is the fastest sketch and is very accurate in all cases, with high precision and good frequency estimation accuracy. CS has low space consumption and is very accurate in most cases, but has slow update rate

¹ Notice that even though we average over 20 runs, we simply partition the input and use the same randomly chosen hash functions throughout the experiment. Hence, failures can propagate over all runs, since we simply query the same sketch repetitively.

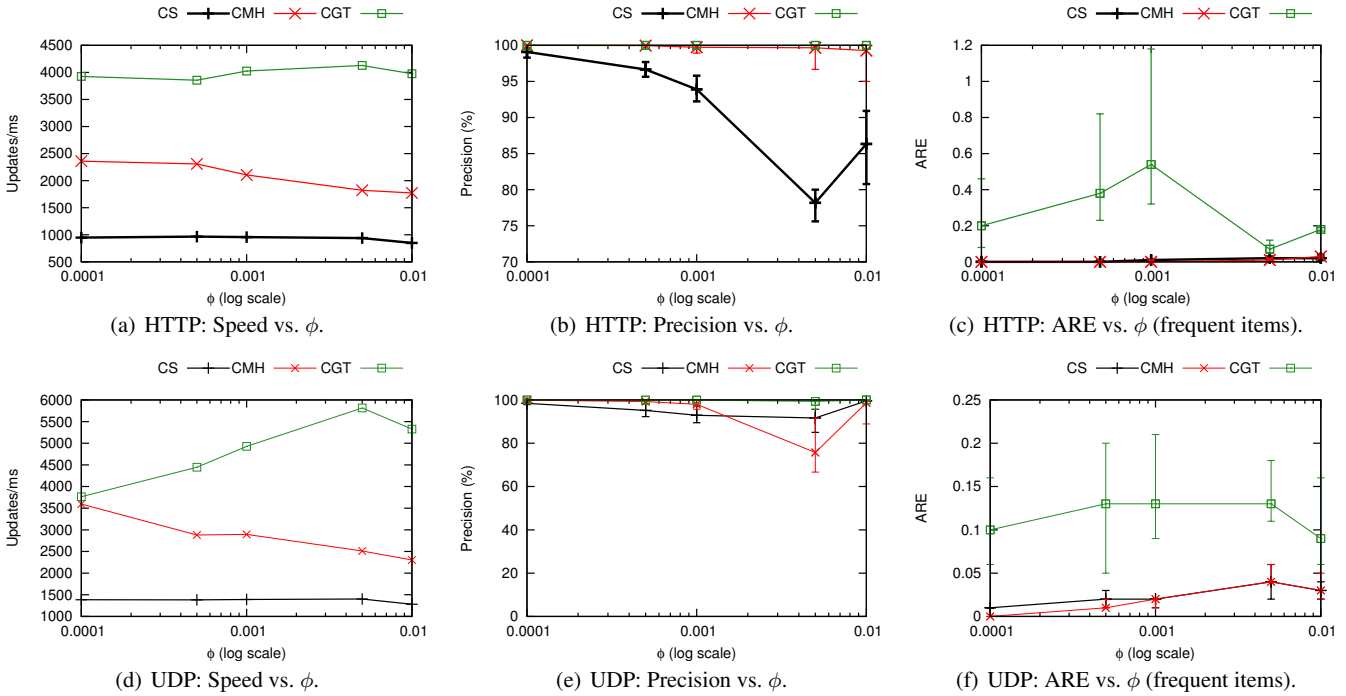


Fig. 14 Performance of sketch algorithms on real network data

and exhibits some random behaviour. The right tradeoff between choice of algorithm, ϕ and hierarchy parameter b will therefore depend on the exact situation.

5 Extensions

As mentioned in the introduction, there are many natural variations of the fundamental frequent items problem which can be addressed by extending the known algorithms.

Weighted input. The definition above assigns all arriving items equal, unit weight. A more general model allows each item to have a weight w , and the frequency of an item is the sum of its weights throughout the stream (and N is replaced with the sum of all weights). The analysis of SPACESAVING and QDIGEST extends to allow arbitrary weights [17]. It is less clear how to generalize LOSSYCOUNTING and FREQUENT for this problem. Sketch methods such as COUNTSKETCH and COUNTMIN directly handle weighted updates: when updating entries in the sketch, the value of the hash function $g_j(i)$ (+1 or -1) is multiplied by the weight w .

Heavy Changes. Given two streams representing, say, the distribution of search queries on two different days, it is possible to subtract one frequency vector from the other to define a *difference vector*. The heavy hitters in this vector correspond to items which have experienced large change in their frequency: this can be positive or negative. This problem has been proposed as a challenge in the context

of web search [27]. Sketch algorithms still give good estimation of frequencies, but the approach of a hierarchical search as advocated by Bu *et al.* [11] may not always work, due to the negative frequencies. Instead, variations of the GROUPTESTING approach can solve this problem [18]; generalizations of this approach trade off the space needed at the cost of increased query time [40].

Top- k items. The top- k items are those items with the k -highest frequencies. Let f_k denote the k th highest frequency; the problem can be restated as finding all items whose frequency exceeds f_{k+1} . But by a similar hardness proof to that for the original exact frequent items problem (Section 2), even approximating this problem by promising to return items whose frequency is at least $(1 - \epsilon)f_k$ requires space linear in the input size. Further weakenings of this problem to make it tractable yields formalizations which are more similar in nature to the approximate frequent items problem.

Assumptions on the frequency distribution. Many realistic frequency distributions are skewed, with a few items with high frequency, and many with low frequency. Such distributions are characterized by the Zipfian, pareto, or power-law distributions (these three distributions are essentially identical up to change of parameters). The Zipfian distribution with parameter z , for example, states that f_k , the k th most frequent item, has frequency proportional to k^{-z} . For large enough z (greater than 1, say), this can simplify the frequent items problem, and reduce the space needed to, for example, $O(\frac{1}{\epsilon^{1/z}})$ in the case of SPACESAVING [35]. Under the

assumption that the frequency distribution is skewed, this can also make the top- k items problem tractable [14, 35].

Tail Bounds. A related issue to top- k and fixed distribution versions is whether distribution-dependent guarantees can be proved. The “residual norm” of a stream, denoted by $F_p^{res(k)}$, is the sum of p th powers of the frequencies after removing the k largest. The guarantees given so far have been in terms of the $F_1^{res(0)}$ norm, but it is straightforward to extend the analysis of the COUNTMIN sketch and the COUNTSKETCH to give bounds in terms of the $F_1^{res(k)}$ and $F_2^{res(k)}$ norms respectively, provided the width of the sketch w is sufficiently larger than k . It is of interest to study whether similar bounds hold for other frequent item algorithms. Such bounds are of use not simply due to the tighter guarantees that follow, but because such tail bounds are at the heart of recent work on signal processing, such as Compressed Sensing [23, 12].

Distributed Streams. A variation of the problem is when there are multiple streams observed by different parties, and the goal is to compute the frequent items over the union of the streams. It is straightforward to solve this problem by combining all the frequent items and their estimated counts from each observer, so an additional requirement is to produce summaries of streams which can be merged together to form summaries of the union of the input streams while occupying no more space than the summary of a single stream. This is relevant when there are many streams and the information needs to be sent over a (sensor) network. Sketches, and the QDIGEST have this merging property. It is less clear how to correctly merge the other counter-based and quantile algorithms while keeping their size bounded.

Distinct Frequent Items. In the Distinct Frequent Items problem (also known as “superspreaders”) the input streams are of the form (i, j) , and f_i is now defined as $|\{j | (i, j) \in S\}|$. Multiple occurrences of (i, j) therefore still only count once towards f_i . Techniques for “distinct frequent items” rely on combining frequent items algorithms with “count distinct” algorithms [31, 7].

Time-decay. While processing a long stream, it may be desirable to weight more recent items more heavily than older ones. Various models of time decay have been proposed. In a sliding window, only the W most recent items, or only the (W) items arriving within the last T time units, should be considered to define the frequent items. The space used should be sublinear in W , and solutions have been proposed with dependency $\log W$ or better [20, 2, 32]. Exponential decay gives an item with ‘age’ a a weight of $\exp(-\lambda a)$ for a fixed parameter λ . The ‘age’ can be derived from timestamps, or implied by the count of items which arrive subsequently. This generates a weighted instance of frequent items, but the weights vary as time increases. However, due

to the structure of the decay function, the decay can be handled quite efficiently [17]. Other decay functions (such as a polynomially decaying weight) require significantly more complex solutions.

6 Conclusions

We have attempted to survey algorithms for finding frequent items in streams, and give an experimental comparison of their behavior to serve as a baseline for comparison. Even so, we had to omit a few less popular algorithms based on random sampling. For insert-only streams, the clear conclusion of our experiments is that the SPACESAVING algorithm, a relative newcomer, has surprisingly clear benefits over others. We observed that implementation choices, such as whether to use a heap or lists of items grouped by frequencies, tradeoff speed and space. Quantile algorithms, with guarantees which appear similar on paper, are demonstrated to be a poor solution for finding frequent items in comparison to the dedicated solutions. They are therefore not recommended for solving this core problem; on the other hand, if a quantile summary is already being kept to summarize a data distribution, it may be used secondarily to extract information about item frequencies. For sketches, there is not such a clear answer, with different algorithms excelling at different aspects of the problem. Due to the increased costs to obtain comparable results, in many situations it seems that counter-based methods are preferable to sketches, unless either the streams contain a large number of deletions (making counter-based methods unviable) or else if hardware issues (such as parallelization or other techniques) make the sketch methods more powerful.

We do not consider this the end of the story, and continue to experiment with other implementation choices. Our source code, datasets and experimental test scripts are available so that others can use these as baseline implementations, and for experimental repeatability. We have done some testing over different computing architectures and observed similar relative performance of the algorithms in terms of throughput.

There have been some other careful comparisons of the performance of streaming algorithms for different problems recently. Dobra and Rusu [22] have studied sketches for the problem of estimating join sizes (vector inner products). Metwally *et al.* compare a variety of algorithms for estimating the number of distinct elements [36]. This recent interest highlights the importance of benchmarking. It is indicative that streaming has “come of age”, in that there are several competing solutions for these fundamental problems, and that these are sufficiently powerful and stable to make it valuable to perform rigorous comparisons. The next logical step is to extend such studies (and availability of code) for other foundational streaming problems, such as finding

quantiles, frequency moments, and more complex mining problems (mining frequent itemsets and clusters).

Acknowledgments. We thank Flip Korn for providing reference implementations of the GK algorithm, and the anonymous reviewers for helpful suggestions.

References

1. N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. In *ACM Symposium on Theory of Computing*, pages 20–29, 1996. Journal version in *Journal of Computer and System Sciences*, 58:137–147, 1999.
2. A. Arasu and G. S. Manku. Approximate counts and quantiles over sliding windows. In *ACM Principles of Database Systems*, 2004.
3. B. Babcock, S. Babu, M. Datar, R. Motwani, and J. Widom. Models and issues in data stream systems. In *ACM Principles of Database Systems*, pages 1–16, 2002.
4. N. Bandi, A. Metwally, D. Agrawal, and A. E. Abbadi. Fast data stream algorithms using associative memories. In *ACM SIGMOD International Conference on Management of Data*, 2007.
5. S. Bhattacharya, A. Madeira, S. Muthukrishnan, and T. Ye. How to scalably skip past streams. In *Scalable Stream Processing Systems (SSPS) Workshop with ICDE 2007*, 2007.
6. L. Bhuvanagiri, S. Ganguly, D. Kesh, and C. Saha. Simpler algorithm for estimating frequency moments of data streams. In *ACM-SIAM Symposium on Discrete Algorithms*, 2006.
7. A. Blum, P. Gibbons, D. Song, and S. Venkataraman. New streaming algorithms for fast detection of superspreaders. Technical Report IRP-TR-04-23, Intel Research, 2004.
8. P. Bose, E. Kranakis, P. Morin, and Y. Tang. Bounds for frequency estimation of packet streams. In *SIROCCO*, 2003.
9. R. S. Boyer and J. Moore. A fast majority vote algorithm. Technical Report ICSCA-CMP-32, Institute for Computer Science, University of Texas, Feb. 1981.
10. R. S. Boyer and J. S. Moore. MJRTY - a fast majority vote algorithm. In *Automated Reasoning: Essays in Honor of Woody Bledsoe*, Automated Reasoning Series, pages 105–117. Kluwer Academic Publishers, 1991.
11. T. Bu, J. Cao, A. Chen, and P. P. C. Lee. A fast and compact method for unveiling significant patterns in high speed networks. In *IEEE INFOCOMM*, 2007.
12. E. Candès and T. Tao. Near optimal signal recovery from random projections and universal encoding strategies. Technical Report math.CA/0410542, arXiv, <http://arxiv.org/abs/math.CA/0410542>, 2004.
13. A. Chakrabarti, G. Cormode, and A. McGregor. A near-optimal algorithm for computing the entropy of a stream. In *ACM-SIAM Symposium on Discrete Algorithms*, 2007.
14. M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. In *Proceedings of the International Colloquium on Automata, Languages and Programming (ICALP)*, 2002.
15. G. Cormode, F. Korn, S. Muthukrishnan, T. Johnson, O. Spatscheck, and D. Srivastava. Holistic UDAFs at streaming speeds. In *ACM SIGMOD International Conference on Management of Data*, pages 35–46, 2004.
16. G. Cormode, F. Korn, S. Muthukrishnan, and D. Srivastava. Space- and time-efficient deterministic algorithms for biased quantiles over data streams. In *ACM Principles of Database Systems*, 2006.
17. G. Cormode, F. Korn, and S. Tirhapura. Exponentially decayed aggregates on data streams. In *IEEE International Conference on Data Engineering*, 2008.
18. G. Cormode and S. Muthukrishnan. What’s new: Finding significant differences in network data streams. In *Proceedings of IEEE Infocom*, 2004.
19. G. Cormode and S. Muthukrishnan. An improved data stream summary: The count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
20. M. Datar, A. Gionis, P. Indyk, and R. Motwani. Maintaining stream statistics over sliding windows. In *ACM-SIAM Symposium on Discrete Algorithms*, 2002.
21. E. Demaine, A. López-Ortiz, and J. I. Munro. Frequency estimation of internet packet streams with limited space. In *European Symposium on Algorithms (ESA)*, 2002.
22. A. Dobra and F. Rusu. Statistical analysis of sketch estimators. In *ACM SIGMOD International Conference on Management of Data*, 2007.
23. D. Donoho. Compressed sensing. <http://www-stat.stanford.edu/~donoho/Reports/2004/CompressedSensing091604.pdf>, 2004. Unpublished Manuscript.
24. M. Fischer and S. Salzberg. Finding a majority among n votes: Solution to problem 81-5. *Journal of Algorithms*, 3(4):376–379, 1982.
25. A. C. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. Strauss. How to summarize the universe: Dynamic maintenance of quantiles. In *International Conference on Very Large Data Bases*, pages 454–465, 2002.
26. M. Greenwald and S. Khanna. Space-efficient online computation of quantile summaries. In *ACM SIGMOD International Conference on Management of Data*, 2001.
27. M. Henzinger. Algorithmic challenges in search engines. *Internet Mathematics*, 1(1):115–126, 2003.
28. J. Hershberger, N. Shrivastava, S. Suri, and C. Toth. Adaptive spatial partitioning for multidimensional data streams. In *ISAAC*, 2004.
29. T. S. Jayram, A. McGregor, S. Muthukrishnan, and E. Vee. Estimating statistical aggregates on probabilistic data streams. In *ACM Principles of Database Systems*, 2007.
30. R. Karp, C. Papadimitriou, and S. Shenker. A simple algorithm for finding frequent elements in sets and bags. *ACM Transactions on Database Systems*, 28:51–55, 2003.
31. G. Kollios, J. Byers, J. Considine, M. Hadjieleftheriou, and F. Li. Robust aggregation in sensor networks. *IEEE Data Engineering Bulletin*, 28(1), Mar. 2005.
32. L. Lee and H. Ting. A simpler and more efficient deterministic scheme for finding frequent items over sliding windows. In *ACM Principles of Database Systems*, 2006.
33. G. Manku and R. Motwani. Approximate frequency counts over data streams. In *International Conference on Very Large Data Bases*, pages 346–357, 2002.
34. G. S. Manku. Frequency counts over data streams. <http://www.cse.ust.hk/vldb2002/VLDB2002-proceedings/slides/S10P03slides.pdf>, 2002.
35. A. Metwally, D. Agrawal, and A. E. Abbadi. Efficient computation of frequent and top-k elements in data streams. In *International Conference on Database Theory*, 2005.
36. A. Metwally, D. Agrawal, and A. E. Abbadi. Why go logarithmic if we can go linear?: Towards effective distinct counting of search traffic. In *International Conference on Extending Database Technology*, 2008.
37. J. Misra and D. Gries. Finding repeated elements. *Science of Computer Programming*, 2:143–152, 1982.
38. S. Muthukrishnan. Data streams: Algorithms and applications. In *ACM-SIAM Symposium on Discrete Algorithms*, 2003.
39. R. Pike, S. Dorward, R. Griesemer, and S. Quinlan. Interpreting the data: Parallel analysis with sawzall. *Dynamic Grids and Worldwide Computing*, 13(4):277–298, 2005.

40. R. Schweller, Z. Li, Y. Chen, Y. Gao, A. Gupta, Y. Zhang, P. A. Dinda, M.-Y. Kao, and G. Memik. Reversible sketches: enabling monitoring and analysis over high-speed data streams. *IEEE Transactions on Networks*, 15(5):1059–1072, 2007.
41. N. Shrivastava, C. Buragohain, D. Agrawal, and S. Suri. Medians and beyond: New aggregation techniques for sensor networks. In *ACM SenSys*, 2004.
42. M. Thorup. Even strongly universal hashing is pretty fast. In *ACM-SIAM Symposium on Discrete Algorithms*, 2000.