

# BRICK: A Novel Exact Active Statistics Counter Architecture

Nan Hua, Jun (Jim) Xu, Bill Lin, and Haiquan (Chuck) Zhao

**Abstract**—In this paper, we present an exact active statistics counter architecture called **Bucketized Rank Indexed Counters (BRICK)** that can efficiently store per-flow variable-width statistics counters entirely in SRAM while supporting both fast updates and lookups (e.g., 40-Gb/s line rates). BRICK exploits statistical multiplexing by randomly bundling counters into small fixed-size buckets and supports dynamic sizing of counters by employing an innovative indexing scheme called rank indexing. Experiments with Internet traces show that our solution can indeed maintain large arrays of exact active statistics counters with moderate amounts of SRAM.

**Index Terms**—Router, statistics counter.

## I. INTRODUCTION

IT IS widely accepted that network measurement is essential for the monitoring and control of large networks. For implementing various network-measurement, router-management, and data-streaming algorithms, there is often a need to maintain very large arrays of statistics counters at wirespeeds (e.g., a million counters for per-flow measurements). For example, on a 40-Gb/s OC-768 link, a new packet can arrive every 8 ns, and the corresponding counter updates need to be completed within this time. While implementing large counter arrays in SRAM can satisfy performance needs, the amount of SRAM required for worst-case counter sizes is often both infeasible and impractical. Therefore, researchers have actively sought alternative ways to realize large arrays of statistics counters at wirespeeds [21]–[23], [28].

In particular, several SRAM-efficient designs of large counter arrays based on hybrid SRAM/DRAM counter architectures have been proposed. Their baseline idea is to store some lower-order bits (e.g., 9 b) of each counter in SRAM, and all its bits (e.g., 64 b) in DRAM. The increments are made only to these SRAM counters, and when the values of SRAM counters become close to overflow, they will be scheduled to be “committed” back to the corresponding DRAM counter.

Manuscript received January 17, 2010; revised September 14, 2010; accepted September 27, 2010; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor T. Wolf. Date of publication March 07, 2011; date of current version June 15, 2011. This work was supported in part by the U.S. National Science Foundation under Grant 0905169. Part of this paper was presented at the ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS), Princeton, NJ, November 6–7, 2009.

N. Hua and J. Xu are with the College of Computing, Georgia Institute of Technology, Atlanta, GA 30332 USA (e-mail: nanhua@cc.gatech.edu; jx@cc.gatech.edu).

B. Lin is with the University of California, San Diego, San Diego, CA 92093 USA (e-mail: billlin@ece.ucsd.edu).

H. Zhao was with the College of Computing, Georgia Institute of Technology, Atlanta, GA 30332 USA. He is now with Microsoft, Redmond, WA 98052 USA.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNET.2011.2111461

These schemes all significantly reduce the SRAM cost. For example, the scheme by Zhao *et al.* [28] achieves the theoretically minimum SRAM cost of between 4–6 b per counter, when the speed difference between SRAM and DRAM ranges between 10 (50 ns/5 ns) and 50 (100 ns/2 ns). However, in these schemes, while writes can be done as fast as on-chip SRAM latencies (2–5 ns), read accesses can only be done as slowly as DRAM latencies (e.g., 60–100 ns). Therefore, such schemes only solve the problem of so-called *passive counters* in which full counter values in general do not need to be read out frequently (not until the end of a measurement epoch). Besides the problem of slow reads, hybrid architectures also suffer from the problem of significantly increasing the amount of traffic between SRAM (usually on-chip) and DRAM (usually off-chip) across the system bus. This may become a serious concern in today’s network processors, where system bus and DRAM bandwidth are already heavily utilized for other packet processing functions [28].

While passive counters are good enough for many network monitoring applications, a number of other applications require the maintenance of *active counters*, in which the values of counters may need to be read out as frequently as they are incremented, typically on a per-packet basis. In many network data-streaming algorithms [4], [8], [14], [15], [26], [27], upon the arrival of each packet, values need to be read out from some counters to decide on actions that need to be taken. For example, if Count-Min sketch [4] is used for elephant detection, we need to read the counter values on a per-packet basis because such readings will decide whether a flow needs to be inserted into a priority queue (implemented as a heap) that stores “candidate elephants.” A prior work on approximate active counters [24] identifies several other data-streaming algorithms that need to maintain active counters, including multistage filters for elephant detection [8] and online hierarchical heavy hitter identification [26]. Currently, all existing algorithms that use active counters implement them as full-size SRAM counters. An efficient solution for exact active counters clearly will save memory cost for all such applications.

## A. Our Approach and Contributions

In this paper, we propose the first solution to the open problem of how to efficiently maintain exact active counters. Our objective is to design an exact counter array scheme that allows for extremely fast read *and* write accesses (at on-chip SRAM speeds). However, these goals will clearly push us back to the origins of using an array of full-size counters in SRAM if we do not impose any additional constraint on the counter values. Fast read access demands that the counters reside entirely in SRAM, and we can make the values of each counter large enough (and random enough) so that each of them needs the worst-case (full-

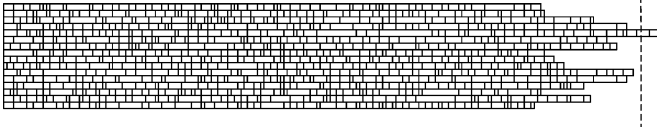


Fig. 1. BRICK wall (conceptual baseline scheme).

size) counter size. Therefore, we will solve our problem under a very natural and reasonable constraint. We assume that the total number of increments, which is exactly the sum of counter values in the array, is bounded by a constant  $M$  during the measurement interval.

This constraint is a reasonable constraint for several reasons. First, this constraint is natural since the number of increments is bounded by the maximum packet arrival rate times the length of the measurement epoch. We can easily enforce an overall count sum limit by limiting the length of the measurement epoch. Moreover, this constraint has been assumed in designing other memory-efficient data structures such as spectral Bloom filters (SBFs) [3]. Furthermore, our scheme will work for arbitrarily large  $M$  values, although its relative memory savings compared to full-size counters get gradually lower with larger  $M$  values.

Let  $N$  be the total number of counters in the array. Then, the ratio  $\frac{M}{N}$  corresponds to the (worst-case) average value of a counter, which is indeed a more relevant parameter than  $M$  for evaluation purposes, as it corresponds to the “per-counter workload.” We observe that small  $\frac{M}{N}$  ratio is dictated by many real-world applications. For example, if we use a Count-Min [4] sketch with  $\ln \frac{1}{\delta}$  arrays of  $\frac{\epsilon}{e}$  ( $e \approx 2.718$ ) counters each for estimating the sizes of TCP/UDP flows, then with probability at least  $1 - \delta$ , the CM-sketch overcounts (it never undercounts) by at most  $M\epsilon$ . Suppose we set  $\delta$  to 0.1 and  $\epsilon$  to  $10^{-5}$  so that we use a total of  $\ln(\frac{1}{0.1}) \times \frac{\epsilon}{10^{-5}} \approx 6.259 \times 10^5$  counters. When the total number of increments  $M$  is set to  $10^8$  and correspondingly the average counts per counter  $\frac{M}{N}$  is approximately 160, we can guarantee that the error is no more than 1000 ( $= 10^8 \times 10^{-5}$ ) with probability at least 0.9. However, 1000 are considered very large errors, and hence for practice we always want  $\frac{M}{N}$  to be much smaller.

We emphasize that even when the ratio  $\frac{M}{N}$  is small, it is still important to figure out ways to save memory, as naive implementations can be grossly wasteful. For example, let the total counts be  $M = 16$  million and the number of counters be  $N = 1$  million. In other words, the average counter value  $\frac{M}{N}$  is 16. Since all increments can go to the same counter, fixed-counter-size design would require a conservative counter size of  $\lg(16 \times 10^6) = 24$  b. However, as we will show, our scheme can significantly reduce the SRAM requirement, which is very important for ASIC implementations where SRAM cost is among the primary costs.

In this paper, we present an exact active counter architecture called Bucketized (B) Rank (R) Indexed (I) Counter (CK), or BRICK. It is built entirely in SRAM so that both read and increment accesses can be processed at tens to hundreds of millions of packets per second. In addition, since it is stored entirely in SRAM, it will not introduce traffic between SRAM and DRAM. This makes it also a very attractive solution for passive counting applications in which the aforementioned problem of increased

traffic over system bus caused by the hybrid SRAM/DRAM architecture becomes a serious concern.

The basic idea of our scheme is intuitive and is based on a very familiar networking concept: statistical multiplexing. Our idea is to bundle groups of a fixed number (let it be 64 in this case) of counters, which is randomly selected from the array, into buckets. We allocate just enough bits to each counter in the sense that if its current value is  $C_i$ , we allocate  $\lfloor \log_2 C_i \rfloor + 1$  b to it. Therefore, counters inside a bucket have variable widths. Suppose the mean width of a counter averaged over the entire array is  $\gamma$ . By the law of large numbers, the total widths of counters in most of the buckets will be fairly close to  $\gamma$  multiplied by the number of counters per bucket. Depicting each counter as a “brick,” as shown in Fig. 1, a section of the “brick wall” illustrates the effect of statistical multiplexing, where each horizontal layer of bricks (consisting of 64 of them) corresponds to a bucket and the length of bricks corresponds to the real counter widths encoding flow sizes in a real-world Internet packet trace (the USC trace in Section V-B).

As we see in this figure, when we set the bucket size to be slightly longer than  $64\gamma$  (the vertical dashed line), the probability of the total widths of the bricks overflowing this line is quite small; among the 20 buckets shown, only one of them has an overflow. Although overflowed buckets need to be handled separately and will cost more memory, we can make this probability small, and the overall overflow cost is small and bounded. Therefore, our memory consumption only needs to be slightly larger than  $64\gamma$  per bucket.

This baseline approach is hard to implement in hardware in practice for two reasons. First, we need to be able to randomly access (i.e., jump to) any counter with ease. Since counters are of variable sizes, we still need to spend several bits per counter for the indexing within the bucket. Note that being able to randomly access is different from being able to delimit all these counters. The latter can be solved with by prefix-free coding (e.g., Huffman coding [5]) of the counter values. Those coding techniques would replace the counter values with variable-length symbols, which could make the size of storage much smaller while making the overhead of accessing and modifying data much larger.

BRICK addresses these two difficulties with a little more overall SRAM cost. It allows for very efficient read and expansion (for increments that increase the width of a counter such as from 15 to 16). A key technique in our data structure is an indexing scheme called *rank indexing*, borrowed from the compression techniques in [7], [12], [13], and [25]. The operations involved in reading and updating this data structure are not only simple for ASIC implementations, but are also supported in modern processors through built-in instructions such as “shift” and “popcount” so that software implementation is efficient (as the involved basic operations such as shift and popcount are supported by modern processors [1], [2]). Therefore, our scheme can be implemented efficiently both in hardware or software.

## B. Background and Related Work

In this section, we compare and contrast our work with previous approaches. One category of approaches is based on the idea of an SRAM/DRAM hybrid architecture [21]–[23], [28].

The state-of-art scheme [28] only requires  $\log_2 \mu$  b per counter, where  $\mu$  is the speed different between SRAM and DRAM. This translates into between 4–6 SRAM bits per SRAM counter. However, the read can take quite long (say at least 100 ns). Therefore, these approaches only solve the passive counting problem.

Another category of approaches is existing active counter solutions [6], [18], [24], which are all based on the approximate counting idea invented by Morris [18]. The idea is to probabilistically increment a counter based on the current counter value. However, approximate counting in general has a very large error margin when the number of bits used is small because the possible estimation values are very sparsely distributed in the range of possible counts. Therefore, when the counter values are small (say 5), its estimation can have a very high relative error (well over 100%). This is not acceptable in network accounting and data-streaming applications where small counter values can be important for overall measurement accuracy. In fact, when the (worst-case) average counter value  $\frac{M}{N}$  is no more than 128, the SRAM cost of our BRICK scheme (about 12 b) is no more than that of [24], which is approximate.

Recently, another counter architecture called counter braids [16] has been proposed, which is inspired by the construction of LDPC codes [9] and can keep track of exact counts of all flows without remembering the association between flows and counters. At each packet arrival, counter increments can be performed quickly by hashing the flow label to several counters and incrementing them. The counter values can be viewed as a linear transformation of flow counts, where the transformation matrix is the result of hashing *all* flow labels during a measurement epoch. However, counter braids are not active and are in fact “more passive” than the SRAM/DRAM hybrid architectures. To find out the size of a single flow, one needs to decode all the flow counts through a fairly long iterative decoding procedure.<sup>1</sup>

Finally, spectral Bloom filter [3] has been proposed, which provides an internal data structure for storing variable-width counters. It uses a hierarchical indexing structure to locate counters that are packed next to each other, which allows for fast random accesses (reads). However, an update that causes the width of the counter  $i$  to grow will cause a shift to counters  $i+1$ ,  $i+2$ , ..., which can have a global cascading effect even with some slack bits provided in between, making it prohibitively expensive when there can be millions of counters. As acknowledged in [3], although the expected amortized cost per update remains constant, and the global cascading effect is small in the average case, the worst case cannot be tightly bounded. Therefore, SBF with variable-width encoding is not an active counter solution as it cannot ensure fast per-packet write accesses at every packet arrival, forcing it to become a mostly read-only data structure in the sense that updates should be orders of magnitude less frequent than queries.

The rest of the paper is organized as follows. Section II describes the design of our scheme in detail. Section III establishes the tail probabilities that allow us to bound and optimize the SRAM requirement. Section IV derives several lower bounds on memory usage for counters to help us understand how far we

are from the optimal. Section V evaluates our scheme by presenting numerical results on memory costs and tail probabilities under various parameter settings, including those extracted from real-world traffic traces.

## II. DESIGN OF BRICK

In this section, we describe the proposed BRICK counter architecture. The objective of BRICK is to efficiently encode a set of  $N$  *exact active* counters  $C_1, C_2, \dots, C_N$ , under the constraint that throughout a network measurement epoch, the total counts<sup>2</sup> across all counters  $\sum_{i=1}^N C_i$  are no more than a predetermined threshold  $M$ , which is carefully justified in Section I. As we explained earlier, since all increments can go to the same counter, the value of a counter can be as large as  $M$ , and hence the worst-case counter width is  $L = \lfloor \log_2 M \rfloor + 1$ . However, it is unnecessarily expensive to allocate  $L$  b to every counter since only a tiny number of them will have counts large enough to require this worst-case width while most others need significantly fewer bits. Therefore, BRICK adopts a sophisticated *variable-width encoding* of counters and can statistically multiplex these variable-width counters through a bucketing scheme to achieve a much more compact representation. However, unlike the aforementioned baseline bucketing scheme, BRICK is extremely SRAM-efficient, yet allows for very fast counter lookup and increment operations.

In the following, we will first present an overview of our proposed design in Section II-A, followed by how it handles lookups, increments, and bucket overflows in Sections II-B–II-D, respectively.

### A. Overview

The basic idea of BRICK is to *randomly* bundle  $N$  counters into  $h$  buckets,  $B_1, B_2, \dots, B_h$ , where each bucket holds  $k$  counters (e.g.,  $k = 64$  in practice) and  $N = hk$ . In each bucket, some counters will be long (possibly  $L$  b in the worst case) and some will be short, depending on the values they contain. As discussed earlier, the objective of bundling is to “statistically multiplex” the variable counter widths in a bucket so that each bucket only needs to be allocated memory space that is slightly larger than  $k$  times the average counter width (across  $N$  counters). Note that since we do not know the actual average width of a counter in advance, we need to instead use the *average width* in the following adversarial context. Imagine that an adversary chooses  $C_1, C_2, \dots, C_N$  values under the constraint  $\sum_{i=1}^N C_i \leq M$  that maximizes the metrics (e.g., average counter width). We emphasize that such an adversary is defined entirely in the well-established context of randomized online algorithm design [19] and has nothing to do with its connotation in security and cryptography.

Fig. 2 depicts these ideas of randomization and bucketization. In particular, as depicted in Fig. 2(a), to access the  $y$ th counter, a pseudorandom permutation function  $\pi : \{1 \dots N\} \rightarrow \{1 \dots N\}$  is first applied to the index  $y$  to obtain a permuted index  $i$ . This pseudorandom permutation function in practice can be as simple<sup>3</sup> as reversing the bits of  $y$ . The corresponding counter  $C_i$

<sup>2</sup>Here, with an abuse of notation, we will use  $C_i$  to denote both the counter and its current count (value).

<sup>3</sup>Since the adversary is defined in the online algorithm context discussed above, we do not believe cryptographically strong pseudorandom permutations, which may increase our cost and slow down our operations, are needed here.

<sup>1</sup>In [16], they need 25 s on a 2.6-GHz computer to decode the flow counts inside a 6-min-long traffic trace.

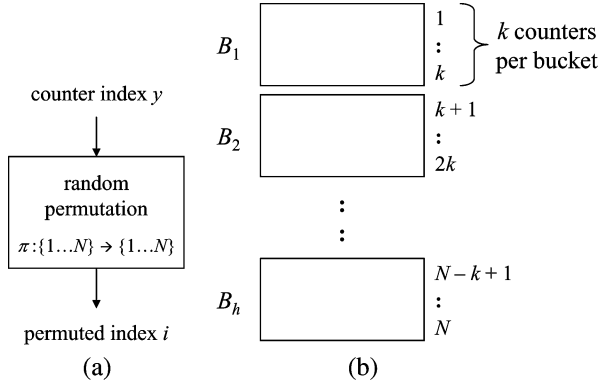


Fig. 2. Randomly bundling counters into buckets. (a) Index permutation. (b) Bucketization.

can then be found in the  $\ell$ th bucket  $B_\ell$ , where  $\ell = \lceil \frac{i}{k} \rceil$ . The bucket structure is depicted in Fig. 2(b). Unless otherwise noted, when we refer to the  $i$ th counter  $C_i$ , we will assume  $i$  is already the result of a random permutation.

As we explained before, the baseline bucketing scheme does not allow for efficient read and write (increment) accesses. In BRICK, a multilevel partitioning scheme is designed to address this problem as follows. The worst-case counter width  $L$  is divided into  $p$  parts, which we refer to as “subcounters.” The  $j$ th subcounter,  $j \in [1, p]$  (from the least significant bits to most significant bits) has  $w_j$  b, such that  $0 < w_j \leq L$  and  $\sum_{j=1}^p w_j = L$ . To save space, for each counter, BRICK maintains just enough of its subcounters to hold its current value. In other words, counters with values no more than  $2^{w_1+w_2+\dots+w_i}$  will not have its  $(i+1)$ th,  $\dots$ ,  $p$ th subcounters stored in BRICK. For example, if  $w_1 = 5$ , any counter with value less than  $2^5 = 32$  will only be allocated a memory entry for its first subcounter. Consider the example shown in Fig. 3(a) with  $k = 8$  counters in a bucket. Only  $C_1$  and  $C_5$  require more than their first subcounters. Such an on-demand allocation requires us to link together all subcounters of a counter, which we achieve using a simple and memory-efficient bitmap indexing scheme called *rank indexing*. Rank indexing enables efficient lookup as well as efficient expansion (when counter values exceed certain thresholds after increments), which will be discussed in detail in Section II-B.

Each bucket contains  $p$  subcounter arrays  $A_1, A_2, \dots, A_p$  to store the first, second,  $\dots$ ,  $p$ th subcounters (as needed) of all  $k$  counters in the bucket. How many entries should be allocated for each array  $A_i$ , denoted as  $k_i$ , turns out to be a nontrivial statistical optimization problem. On the one hand, to save memory, we would like to make  $k_2, k_3, \dots, k_p$  ( $k_1$  is fixed as  $k$ ) as small as possible. On the other hand, when we encounter the unlucky situation that we need to exceed any of these limits (say for a certain  $d$ , we have more than  $k_d$  counters in a bucket that have values larger than or equal to  $2^{w_1+w_2+\dots+w_{i-1}}$ ), then we will have a “bucket overflow” that would require that all counters inside the bucket be relocated to an additional array of full-size buckets with fixed worst-case width  $L$  for each counter, as we will show in Section II-D. Given the high cost of storing a duplicated bucket in the full-size array, we would like to choose larger  $k_2, \dots, k_p$  to make this probability as small as possible. In Section III, we develop extremely tight tail bounds on the overflow probability that allows us to choose parameters  $\{k_i\}_{2 \leq i \leq p}$

and  $\{w_i\}_{1 \leq i \leq p-1}$  to achieve near-optimal tradeoffs between these two conflicting issues and minimize the overall memory consumption.

## B. Rank Indexing

A key technique in our data structure is an indexing scheme that allows us to efficiently identify the locations of the subcounters across the different subcounter arrays for some counter  $C_i$ . In particular, for  $C_i$ , its  $d$  subcounters  $C_{i,1}, \dots, C_{i,d}$  are spread across  $A_1, \dots, A_d$  at locations  $a_{i,1}, \dots, a_{i,d}$ , respectively (i.e.,  $C_{i,j} = A_j[a_{i,j}]$ ). For example, as shown in Fig. 3(b),  $C_5$  is spread across  $A_3[1] = 10$ ,  $A_2[2] = 11$ , and  $A_1[5] = 11011$ .

For each bucket, we maintain an index bitmap  $I$ .  $I$  is divided into  $p-1$  parts,  $I_1, \dots, I_{p-1}$ , with a one-to-one correspondence to the subcounter arrays  $A_1, \dots, A_{p-1}$ , respectively. Each part  $I_j$  is a bitmap with  $k_j$  b,  $I_j[1], \dots, I_j[k_j]$ , one bit  $I_j[a]$  for each entry  $A_j[a]$  in  $A_j$ . Each  $I_j[a]$  is used to determine if the counter stored in  $A_j[a]$  has expanded beyond the  $j$ th subcounter array.  $I_j$  is also used to compute the index location of  $C_i$  in the next subcounter array  $A_{j+1}$ . Because a counter cannot expand beyond the last subcounter array, there is no need for an index bitmap component for the most significant subcounter array  $A_p$ . For example, consider the entries  $A_1[1]$  and  $A_1[5]$ , where the corresponding counter has expanded beyond  $A_1$ . This is indicated by having the corresponding bit positions  $I_1[1]$  and  $I_1[5]$  set to 1, as shown in shaded boxes in Fig. 3(b). All remaining bit positions in  $I_1$  are set to 0, as shown in clear boxes.

For each counter that has expanded beyond  $A_1$ , an arrow is shown in Fig. 3(b) that links a subcounter in  $A_1$  with the corresponding subcounter entry in  $A_2$ . For example, for  $C_5$ , its subcounter entry  $A_1[5]$  in  $A_1$  is linked to the subcounter entry  $A_2[2]$  in  $A_2$ . Rather than expending memory to store these links explicitly, which could vanish savings gained by reduced counter widths, we *dynamically* compute the location of a subcounter in the next subcounter array  $A_{j+1}$  based on the current bitmap  $I_j$ . This way, no memory space is needed to store link pointers. This dynamic computation can be readily determined using an operation called *rank*( $s, j$ ), which returns the number of ones only in the range  $s[1] \dots s[j]$  in the bit-string  $s$ . This operation is similar to the *rank* operator defined in [13].

We apply the *rank* operator on a bitmap  $I_j$  by interpreting it as a bit-string. As we shall see in Sections III and V, our approach is designed to work with small buckets of counters (e.g.,  $k = 64$ ). Therefore, the corresponding bit-strings  $I_j$  are also relatively short since all subcounter arrays satisfy  $k_j \leq k$ . Moreover, each successive  $k_j$  in the higher subcounter arrays is substantially smaller than the previous subcounter array, with the corresponding reduction in the length of the bit-string  $I_j$ . In turn, the *rank* operator can be efficiently implemented by combining a bitwise-AND instruction with another operation called *popcount*( $s$ ), which returns the number of ones in the bit-string  $s$ . Fortunately, the *popcount* operator is becoming an increasingly available hardware-optimized instruction in modern microprocessors and network processors. For example, current generations of 64-b x86 processors have this instruction built in [1], [2]. Using this *popcount* instruction, the *rank* operation for bit-strings with lengths up to  $|s| = 64$  b can be readily computed in as few as two instructions. As shown with numerical examples and trace simulations in Section V, very good results can be achieved with a bucket size fixed at 64.

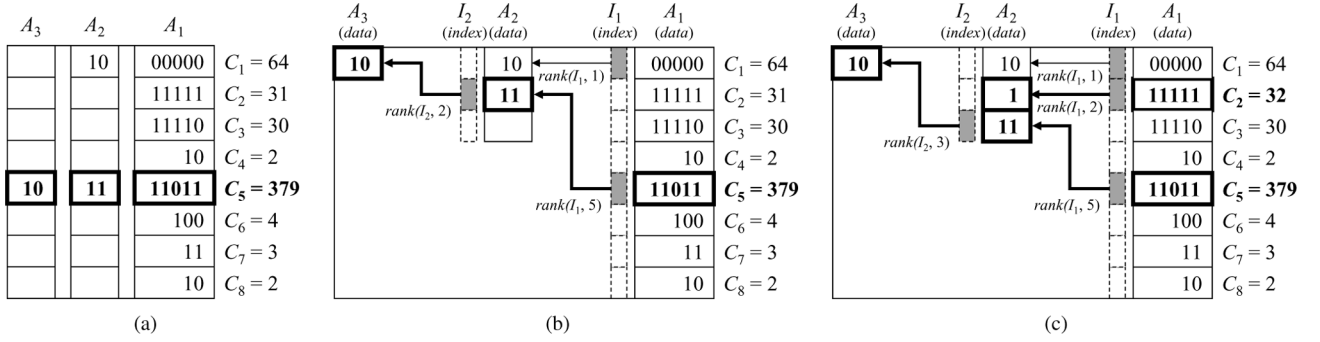


Fig. 3. (a) Within a bucket, segmentation of variable-width counters into subcounter arrays. (b) Compact representation of variable-width counters. (c) Updated data structure after incrementing  $C_2$ .

The pseudocode for the lookup operation is shown in Algorithm 1. The retrieval of the subcounters using rank indexing is shown in lines 3–6, with the final count returned at the end of the procedure. For a hardware implementation, the iterative procedure can be readily pipelined. As we shall see in Section V, we only need a small number of levels (e.g., three) in practice to achieve efficient results.

#### Algorithm 1: Pseudocode

```

1  lookup( $i$ )
2     $C_i = 0$ ;  $a = i \bmod k$ ;
3    for  $j = 1$  to  $p$ 
4       $C_{i,j} = A_j[a]$ ;
5      if ( $j == p$  or  $I_j[a] == 0$ ) break;
6       $a = \text{rank}(I_j, a)$ ;
7    return  $C_i$ ;
8  increment( $i$ )
9     $a = i \bmod k$ ;
10   for  $j = 1$  to  $p$ 
11      $A_j[a] = A_j[a] + 1$ ;
12     if ( $j == p$  or  $A_j[a] \neq 0$ ) break; /* last array
13     or no carry */
14     if ( $I_j[a] == 1$ ) /* next level already allocated */
15        $a = \text{rank}(I_j, a)$ ;
16     else /* expand */
17        $I_j[a] = 1$ ;
18        $a = \text{rank}(I_j, a)$ ;
19        $b = (a - 1)w_{j+1} + 1$ ;
20        $A_{j+1} = \text{varshift}(A_{j+1}, b, w_{j+1})$ ;
21        $I_{j+1} = \text{varshift}(I_{j+1}, a, 1)$ ;
22        $A_{j+1}[a] = 1$ ;
23     break;
```

#### C. Handling Increments

The increment operation is also based on the traversal of subcounters using rank indexing. We will first describe the basic idea by means of an example. Consider the counter  $C_2$  in Fig. 3(b). Its count is 31, which can be encoded in just the subcounter array  $A_1$  with  $C_{2,1} = 11111$ . Suppose we want to increment  $C_2$ . We first increment its first subcounter component  $C_{2,1} = 11111$ , which results in  $C_{2,1} = 00000$  with a carry propagation to the next level. This is depicted in Fig. 3(c).

This carry propagation triggers the increment of the next sub-counter component  $C_{2,2}$ . The location of  $C_{2,2}$  can be determined using rank indexing (i.e.,  $\text{rank}(I_1, 2) = 2$ ). However, the location of  $A_2[2]$  was previously occupied by the counter  $C_5$ . To maintain *rank ordering*, we have to *shift* the entries in  $A_2$  down by one to free up the location  $A_2[2]$ . This is achieved by applying an operation called *varshift*( $s, j, c$ ), which performs a right shift on the substring starting at bit position  $j$  by  $c$  b (with vacant bits filled by zeros). The *varshift* operator can be readily implemented in most processors by means of shift and bitwise-logical instructions.

In particular, we can view a subcounter array  $A_j$  as a *bit-string* formed by the concatenation of its entries, namely  $A_j = A_j[1]A_j[2] \dots A_j[k_j]$ . The starting bit position for an entry  $A_j[a]$  in the bit-string can be computed as  $b = (a - 1)w_j + 1$ , where  $w_j$  is the bit width of the subcounter array  $A_j$ . Consider  $C_5$  in Fig. 3(c). After the shifting operation has been applied, the location of its subcount in  $A_2$  will be shifted down by one entry. Therefore, its corresponding expansion status in  $I_2$  must be shifted down by one position as well. The carry propagation of  $C_2$  into  $A_2$  is achieved by setting  $A_2[2] = 1$ .

As with the rank operator, BRICK has been designed to work with small fixed-size buckets so that *varshift* can be directly implemented using hardware-optimized instructions. In particular, *varshift* only has to operate on  $A_2$  or higher. Since the size of each level decreases exponentially, the bit-strings formed by each subcounter array  $A_2$  and above are also very short. As the results show in Section V, with a bucket size of 64, all sub-counter arrays  $A_2$  and above have a string length at most 64 b, much less for the higher levels. Therefore, *varshift* can be directly implemented using 64-b instructions.

The pseudocode for the increment operation is shown in the latter part of Algorithm 1. Again, the iterative procedure shown in Algorithm 1 for increment is readily amenable to pipelining in hardware. In general, the lookup or update of each successive level of subcounter arrays can be pipelined such that at each packet arrival, a lookup or update can operate on  $A_1$  while a previous operation operates on  $A_2$ , and so forth.

#### D. Handling Overflows

Thus far, we have assumed in our basic data structure that we are guaranteed that each subcounter array has been dimensioned to always provide sufficient entries to store all subcounters in a

bucket. To achieve greater memory efficiency, the number of entries in the subcounter arrays can be reduced so that there is only a very small probability that a bucket will not have sufficient subcounter array entries. As rigorously analyzed in Section III and numerically evaluated in Section V, this bucket overflow probability can be made arbitrarily small while achieving significant reduction in storage for each bucket.

To facilitate this overflow handling, we extend the basic data structure described in Section II-A with a small number of *full-size* buckets  $F_1, F_2, \dots, F_J$ . Each full-size bucket  $F_t$  is organized as  $k$  full-size counters (i.e., all counters with a worst-case width of  $L$  b). When a bucket overflow occurs for some  $B_\ell$ , the next available full-size bucket  $F_t$  is allocated to store its  $k$  counters, where  $t$  is just +1 of the last allocated full-size bucket. An overflow status flag  $f_\ell$  is set to indicate the bucket has overflowed. The index of the full-size bucket  $F_t$  is stored in a field labeled  $t_\ell$ , which is associated with  $B_\ell$ . In practice, we only need a small number of full-size buckets. As shown in Section V, for real Internet traces with over a million counters, only about  $J \approx 100$  full-size buckets are enough to handle the overflow cases. Therefore, the index field only requires a small number of extra bits per bucket (e.g., 7 b).

Rather than migrating *all*  $k$  counters from  $B_\ell$  to  $F_{t_\ell}$  *at once*, a counter is only migrated *on demand* upon the next increment operation (“migrate-on-write”). This way, the migration of an overflow counter to a full-size counter does not disrupt other counter updates. The location of counter  $C_i$  in  $F_{t_\ell}$  is simply  $a = i \bmod k$ , as before. To indicate if counter  $C_i$  has been migrated, a migration status flag  $g_{t_\ell}[a]$  is associated with each counter entry  $F_{t_\ell}[a]$  (i.e.,  $g_{t_\ell}[a] = 1$  indicates that the corresponding counter has been migrated).

The modified lookup operation simply first checks if a counter from an overflowed bucket has already been migrated, in which case the full-size count is simply retrieved from the corresponding full-size bucket entry. Otherwise, the counter is retrieved as before. The modified increment operation is extended in a similar manner. It first checks if a counter from an overflowed bucket has already been migrated, in which case the full-size counter in the corresponding full-size bucket is incremented. If the counter is from a previously overflowed bucket  $B_\ell$ , but it has not been migrated yet, then it is read from  $B_\ell$ , incremented, and migrated-on-write to the corresponding location in the full-size bucket. Otherwise, the counter in  $B_\ell$  is incremented as before. Finally, before propagating a carry to the next level, we first check if all entries in the next subcounter array are already being used. If so, the next full-size bucket is allocated, and the incremented count is migrated-on-write to the corresponding location.

### III. ANALYSIS

#### A. Analytical Guarantees

In this section, we bound the failure probability  $P_f$  that the number of overflowed buckets, each of which carries the hefty penalty of having to be allocated an additional bucket of full-size counters (as discussed in Section II-D), will exceed any given threshold  $J$ . We will establish a rigorous relationship between  $P_f$  and parameters  $k_2, k_3, \dots, k_p$ , the number of entries BRICK allocates to subcounter arrays  $A_2, \dots, A_p$  (the size of  $A_1$  is already fixed to  $k$ ), and  $w_1, w_2, \dots, w_p$ , the widths of an entry in  $A_2, \dots, A_p$ . The ultimate objective of this analysis is to

find the optimal tradeoff between  $k_2, k_3, \dots, k_p$  and  $J$  that allows us to minimize the amount of overall memory consumption ( $h = N/k$  regular buckets +  $J$  full-size buckets) while keeping the failure probability  $P_f$  under an acceptable threshold (say  $10^{-10}$  or even smaller). Surprisingly, the theory of stochastic ordering [20], which seems unrelated to the context of this paper, plays a major role in these derivations.

Recall that the maximum counter width  $L$  is partitioned into subcounter widths  $w_1, w_2, \dots, w_p$ . Only counters whose value is larger than or equal to  $2^{L_d}$ , where  $L_d$  is defined as  $\sum_{j=1}^{d-1} w_j$ , will need an entry in the subcounter array  $A_d$  of a bucket. Since the aggregate count of all counters is no more than  $M$ , we know that there will be at most  $m_d$  of such counters in the whole counter array, where  $m_d$  is defined as  $M2^{-L_d}$ .

Now, imagine at most  $m_d$  such counters are uniformly randomly distributed into  $N$  array locations through the aforementioned index permutation scheme. We hope that they are very evenly distributed among these buckets so that very few buckets will have more than  $k_d$  of them falling into it (i.e., overflow of  $A_d$ ). Suppose we dimension  $J_d$  full-size buckets to handle bucket overflows caused by these counters. We would like to bound the probability that more than  $J_d$  buckets will have their  $A_d$  arrays overflowed.

We will consider the worst-case scenario that there are exactly  $m_d$  counters needing entries in  $A_d$ . If there are less such counters, the overflow probability will only be smaller, and our tail bound still applies. For convenience, we denote the percentage of them in the counter array  $m_d/N$  as  $\alpha_d$ .

Let random variables  $X_{1,d}, X_{2,d}, \dots, X_{h,d}$  be the number of used entries in the subcounter array  $A_d$  among the buckets  $B_1, B_2, \dots, B_h$ . Each array location has a probability  $\alpha_d$  of being assigned one of the  $m_d$  counters, and there are  $k$  array locations in each bucket, so  $X_{j,d}$  is roughly distributed as  $\text{Binomial}(k, \alpha_d)$  for any  $j$ . Here,  $\text{Binomial}(\mathcal{N}, \mathcal{P})$  is the binomial distribution with  $\mathcal{N}$  trials and  $\mathcal{P}$  as the success probability of each trial. Therefore, the overflow probability of level  $d$  from any bucket  $B_j$  is roughly

$$\epsilon_d = \text{Binotail}_{k, \alpha_d}(k_d)$$

where  $\text{Binotail}_{\mathcal{N}, \mathcal{P}}(\mathcal{K}) \equiv \sum_{z=\mathcal{K}+1}^{\mathcal{N}} \binom{\mathcal{N}}{z} \mathcal{P}^z (1-\mathcal{P})^{(\mathcal{N}-z)}$  denotes the tail probability  $\Pr[Z > \mathcal{K}]$ , where  $Z$  has distribution  $\text{Binomial}(\mathcal{N}, \mathcal{P})$ .

Intuitively, these random variables are *almost independent*, as the only dependence among them seems to be that their total is  $m_d$ . If we do assume that they are independent, then the probability that the number of total overflows be larger than  $J_d$  entries is roughly

$$\delta_d = \text{Binotail}_{h, \epsilon_d}(J_d).$$

Readers understandably will immediately protest this voodoo tail bound result since the  $X_{j,d}$ 's are not exactly binomial, and they are not actually independent. Interestingly, we are able to establish a rigorous tail bound of  $2\delta_d$ , which is only two times the voodoo tail bound  $\delta_d$ . A similar bound has been established by Mitzenmacher and Upfal in their book [17], which used independent Poisson distributions to bound multinomial distributions, using techniques from stochastic ordering theory [20] implicitly (i.e., without introducing such concepts). In our case, we use independent binomial distributions to

bound multivariate hypergeometric distributions, i.e., those of  $X_{1,d}, X_{2,d}, \dots, X_{h,d}$ .

Based on this rigorous tail bound to be proven in Section III-B and taking union of the overflow events from all the subarrays, we arrive at the following corollary.

*Corollary 1:* Let parameters  $\delta_2, \dots, \delta_p$  be defined as above. The failure probability of insufficient full-size buckets, i.e., that the total number of overflows that need to be moved to the additional full-size buckets from all subarrays exceeds  $J = J_2 + \dots + J_p$ , is no more than  $2(\delta_2 + \dots + \delta_p)$ .

If given a target worst-case failure probability  $P_f$  of insufficient full-size buckets, e.g.,  $10^{-10}$  or even smaller, an optimization procedure remains to configure parameters from  $w_1$  to  $w_{p-1}$ ,  $k_2$  to  $k_p$ , and  $J_2$  to  $J_p$ , so that we can achieve the best tradeoff for the overall memory space, which takes into consideration the storage of all subcounter arrays, index bitmaps, and all full-size buckets, and even the  $\lceil \lg(J) \rceil + 2$  b for  $f_\ell$  and  $t_\ell$  in each bucket, which indicate the migration to full-size buckets.

Given the messy nature of the Binomial distribution, “clean” analytical solutions (e.g., based on Lagrange multipliers) do not exist. We designed a quick search strategy that can generate near-optimal configurations. Our evaluation results in Section V are obtained based on the near-optimal parameter configurations generated by this procedure. We omit the detail of this procedure in the interest of space.

### B. Our Main Tail Bound

In this section, we formally state the aforementioned tail bound theorem (two times the voodoo bound). We would like to state this theorem using generic parameters that have the same symbol as before, but without the subscript  $d$  since they can be replaced by the corresponding parameters with subscript  $d$  to obtain the tail bound on the number of overflows from every subarray  $A_d$ . In particular, we will replace  $m_d$  (the number of counters that will have an entry in subcounter array  $A_d$ ) by  $m$ , and  $k_d$  (the number of entries in subcounter array  $A_d$ ) by  $c$ , as  $k$  has been used to denote the number of counters in each bucket in the original counter array. Furthermore, to highlight the general nature of our theorem, we further detach ourselves from the application semantics by stating the theorem as follows.

*Theorem 1:*  $m$  balls are uniformly randomly thrown into  $h$  buckets that have  $k$  entries each, with at most one ball in each entry. Let  $N = hk$ . Let  $X_1^{(m)}, \dots, X_h^{(m)}$  be the number of balls that fall into each bucket. Let  $\alpha = (m/hk)$  and assume  $\alpha \leq (1/2)$ . Let  $Y_1^{(\alpha)}, \dots, Y_h^{(\alpha)}$  be independent random variables distributed as Binomial( $k, \alpha$ ). Let  $f(x_1, \dots, x_h)$  be an increasing function in each argument. Then

$$E[f(X_1^{(m)}, \dots, X_h^{(m)})] \leq 2E[f(Y_1^{(\alpha)}, \dots, Y_h^{(\alpha)})].$$

Before we prove this theorem, we need to formally characterize the underlying probability model and in particular specify precisely what we mean by throwing  $m$  balls “uniformly randomly” into  $N$  entries as follows. Among all  $\binom{N}{m}$  ways of injective mapping from  $m$  balls into  $N$  entries, every way happens with equal probability  $\frac{1}{\binom{N}{m}}$ , when these balls are considered indistinguishable. We refer to this characterization of the underlying probability model as “throwing  $m$  balls into  $N$  entries in one shot.” It is not hard to verify that the following process of “throwing  $m$  balls into  $N$  entries one by one” results in the

same probability model. In this process, at first a ball is thrown into an entry chosen uniformly from these  $N$  entries. Then, another ball is thrown into an entry uniformly picked from the remaining  $N - 1$  entries, and so on. This equivalent characterization of the underlying probability model makes it easier for us to establish the stochastic ordering relationship among vectors of random variables in Section III-C, an essential step for the proof of Theorem 1. Now, we are ready to prove Theorem 1.

*Proof of Theorem 1:* We use  $X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}$  when there are  $l$  balls thrown instead of  $m$ . In Proposition 1, we prove that for any  $l$  value,  $\mu(X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)})$  is equivalent to  $\mu(Y_1^{(\alpha)}, Y_2^{(\alpha)}, \dots, Y_h^{(\alpha)} | \sum_{j=1}^h Y_j^{(\alpha)} = l)$ , where  $\mu(Z)$  denotes the distribution of a random variable or vector  $Z$ . In other words, conditioned upon  $\sum_{j=1}^h Y_j^{(\alpha)} = l$ , the independent random variables  $Y_1^{(\alpha)}, Y_2^{(\alpha)}, \dots, Y_h^{(\alpha)}$  have the same joint distribution as dependent random variables  $X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}$ . Then, we prove in Proposition 3 that when  $l \leq l'$ ,  $[X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}]$  is stochastically less than or equal to (defined later)  $[X_1^{(l')}, X_2^{(l')}, \dots, X_h^{(l')}]$ . For any increasing function  $f(x_1, x_2, \dots, x_h)$ , we have

$$\begin{aligned} E[f(Y_1^{(\alpha)}, \dots, Y_h^{(\alpha)})] &= \sum_{l=0}^N E\left[f(Y_1^{(\alpha)}, \dots, Y_h^{(\alpha)}) | \sum_{j=1}^h Y_j^{(\alpha)} = l\right] \Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} = l\right] \\ &\geq \sum_{l=m}^N E\left[f(Y_1^{(\alpha)}, \dots, Y_h^{(\alpha)}) | \sum_{j=1}^h Y_j^{(\alpha)} = l\right] \Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} = l\right] \\ &= \sum_{l=m}^N E[f(X_1^{(l)}, \dots, X_h^{(l)})] \Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} = l\right] \quad (1) \end{aligned}$$

$$\begin{aligned} &\geq \sum_{l=m}^N E[f(X_1^{(m)}, \dots, X_h^{(m)})] \Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} = l\right] \quad (2) \\ &= E[f(X_1^{(m)}, \dots, X_h^{(m)})] \Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} \geq m\right] \\ &= E[f(X_1^{(m)}, \dots, X_h^{(m)})] \text{Binotail}_{N,\alpha}(m-1) \\ &\geq \frac{1}{2} E[f(X_1^{(m)}, \dots, X_h^{(m)})]. \quad (3) \end{aligned}$$

Equation (1) is due to Proposition 1, inequality (2) is due to Proposition 3, and inequality (3) is due to the properties of the 50-percentile point of binomial distributions proven in [10]. ■

*Corollary 2:* Let the variable be as defined in Theorem 1. Let  $c$  and  $J$  be some constants. Let  $\epsilon = \text{Binotail}_{k,\alpha}(c)$ . Then

$$\Pr\left[\sum_{j=1}^h 1_{\{X_j^{(m)} > c\}} > J\right] \leq 2\text{Binotail}_{h,\epsilon}(J).$$

*Proof:* Consider function  $f(x_1, x_2, \dots, x_h) \equiv 1_{\{\sum_{j=1}^h 1_{\{x_j > c\}} > J\}}$ , which is an increasing function of  $x_1, \dots, x_h$ . From Theorem 1, we have  $\Pr\left[\sum_{j=1}^h 1_{\{X_j^{(m)} > c\}} > J\right] \leq 2\Pr\left[\sum_{j=1}^h 1_{\{Y_j^{(\alpha)} > c\}} > J\right]$ .

Since  $\left\{1_{\{Y_j^{(\alpha)} > c\}}\right\}_{1 \leq j \leq h}$  are independent Bernoulli random variables with probability  $\epsilon = \text{Binotail}_{k,\alpha}(c)$ , their sum is distributed as  $\text{Binomial}(h, \epsilon)$ . Therefore,  $\Pr\left[\sum_{j=1}^h 1_{\{Y_j^{(\alpha)} > c\}} > J\right]$  is equal to  $\text{Binotail}_{h,\epsilon}(J)$ . ■

### C. Proofs of Propositions 1–3

**Proposition 1:**  $\mu(X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}) = \mu(Y_1^{(\alpha)}, Y_2^{(\alpha)}, \dots, Y_h^{(\alpha)} | \sum_{j=1}^h Y_j^{(\alpha)} = l)$ .

*Proof:* It suffices to prove that for any nonnegative integers  $l_1, l_2, \dots, l_h$  that satisfy  $\sum_{j=1}^h l_j = l$ ,  $\Pr[X_1^{(l)} = l_1, X_2^{(l)} = l_2, \dots, X_h^{(l)} = l_h] = \Pr[Y_1^{(\alpha)} = l_1, Y_2^{(\alpha)} = l_2, \dots, Y_h^{(\alpha)} = l_h | \sum_{j=1}^h Y_j^{(\alpha)} = l]$ . We show that both the left-hand side (LHS) and the right-hand side (RHS) are equal to

$$\frac{\binom{k}{l_1} \binom{k}{l_2} \cdots \binom{k}{l_h}}{\binom{N}{l}}. \quad (4)$$

Since there are  $\binom{N}{l}$  ways of selecting  $l$  entries out of a total of  $N$  entries, and each way happens with equal probability  $\frac{1}{\binom{N}{l}}$ , the LHS is equal to (4) because there are  $\binom{k}{l_1} \binom{k}{l_2} \cdots \binom{k}{l_h}$  ways among them that result in the event  $\{X_1^{(l)} = l_1, X_2^{(l)} = l_2, \dots, X_h^{(l)} = l_h\}$ . Now, we prove that the RHS is equal to (4) as well. Since  $Y_1^{(\alpha)}, Y_2^{(\alpha)}, \dots, Y_h^{(\alpha)}$  are independent random variables with distribution  $\text{Binomial}(k, \alpha)$ ,  $\sum_{j=1}^h Y_j^{(\alpha)}$  has distribution  $\text{Binomial}(N, \alpha)$  and therefore

$$\Pr\left[\sum_{j=1}^h Y_j^{(\alpha)} = l\right] = \binom{N}{l} \alpha^l (1-\alpha)^{N-l}. \quad (5)$$

Additionally, when  $\sum_{j=1}^h Y_j^{(\alpha)} = l$ , we have

$$\begin{aligned} \Pr\left[Y_1^{(\alpha)} = l_1, Y_2^{(\alpha)} = l_2, \dots, Y_h^{(\alpha)} = l_h, \sum_{j=1}^h Y_j^{(\alpha)} = l\right] \\ = \prod_{j=1}^h \binom{k}{l_j} \alpha^{l_j} (1-\alpha)^{k-l_j} = \alpha^l (1-\alpha)^{N-l} \prod_{j=1}^h \binom{k}{l_j}. \end{aligned} \quad (6)$$

Combining (5) and (6), we obtain that the RHS is equal to (4) as well. ■

Stochastic ordering is a way to compare two random variables. Random variable  $X$  is stochastically less than or equal to random variable  $Y$ , written  $X \leq_{\text{st}} Y$ , iff  $E\phi(X) \leq E\phi(Y)$  for all increasing functions  $\phi$  such that the expectations exists. An equivalent definition of  $X \leq_{\text{st}} Y$  is that  $\Pr[X > t] \leq \Pr[Y > t]$ ,  $-\infty < t < \infty$ . The definition involving increasing functions also applies to random vectors  $X = (X_1, \dots, X_h)$  and  $Y = (Y_1, \dots, Y_h) : X \leq_{\text{st}} Y$  iff  $E\phi(X) \leq E\phi(Y)$  for all increasing functions  $\phi$  such that the expectations exists. Here,  $\phi$  is increasing means that it is increasing in each argument separately with other arguments being fixed. This is equivalent to

$\phi(X) \leq_{\text{st}} \phi(Y)$ . Note this definition is a much stronger condition than  $\Pr[X_1 > t_1, \dots, X_h > t_h] \leq \Pr[Y_1 > t_1, \dots, Y_h > t_h]$  for all  $t = (t_1, \dots, t_h) \in \mathcal{R}^n$ .

Now, we state without proof a fact that will be used to prove Proposition 3. Its proof can be found in all books that deal with stochastic ordering [20].

**Proposition 2:** Let  $X$  and  $Y$  be two random variables (or vectors).  $X \leq_{\text{st}} Y$  iff there exists  $X'$  and  $Y'$  such that  $\mu(X') = \mu(X)$ ,  $\mu(Y') = \mu(Y)$ , and  $\Pr[X' \leq Y'] = 1$ .

Now, we are ready to prove the following proposition.

**Proposition 3:** For any  $0 \leq l < l' \leq N$ , we have

$$[X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}] \leq_{\text{st}} [X_1^{(l')}, X_2^{(l')}, \dots, X_h^{(l')}].$$

*Proof:* It suffices to prove it for  $l' = l + 1$ . Our idea is to find random variables  $Z$  and  $W$  such that  $Z$  has the same distribution as  $[X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}]$ ,  $W$  has the same distribution as  $[X_1^{(l+1)}, X_2^{(l+1)}, \dots, X_h^{(l+1)}]$ , and  $\Pr[Z \leq W] = 1$ . We will use the aforementioned probability model that is generated by “throwing  $m$  balls into  $N$  entries one-by-one” random process. Now, given any outcome  $\omega$  in the probability space  $\Omega$ , let  $Z(\omega) = [Z_1(\omega), Z_2(\omega), \dots, Z_h(\omega)]$ , where  $Z_j(\omega)$  is the number of balls in the  $j$ th bucket after we throw  $l$  balls into these  $N$  entries one by one. Now with all these  $l$  balls there, we throw the  $(l+1)$ th ball uniformly randomly into one of the remaining empty entries. We define  $W(\omega)$  as  $[W_1(\omega), W_2(\omega), \dots, W_h(\omega)]$ , where  $W_j(\omega)$  is the number of balls in the  $j$ th bucket after we throw in the  $(l+1)$ th ball. Clearly, we have  $Z(\omega) \leq W(\omega)$  for any  $\omega \in \Omega$ , and therefore  $\Pr[Z \leq W] = 1$ . Finally, we know from the property of the “throwing  $m$  balls into  $N$  entries one-by-one” random process that  $Z$  and  $W$  have the same distribution as  $[X_1^{(l)}, X_2^{(l)}, \dots, X_h^{(l)}]$  and  $[X_1^{(l+1)}, X_2^{(l+1)}, \dots, X_h^{(l+1)}]$ , respectively. ■

## IV. INFORMATION THEORY BOUND

We are interested in how far we are from the optimal memory usage cost. In this section, we will try to answer this question partly by deriving several lower bounds on the minimum memory requirement per counter under the aforementioned constraint that the sum of counter values  $C_1, C_2, \dots, C_N$  (i.e., the total number of increments) is no more than  $M$ .

We will explore this question in the following sequence. First, we will present a naive bound, which is the worst-case total number of bits needed to store all these counter values. However, this bound does not take into account the indexing cost. Therefore, secondly, we will explore the additional indexing cost based on information theory. However, this part is only aimed at getting a feeling about how much additional indexing cost is required since the derivation is not strict. Finally, we will analyze and prove the minimum number of bits needed to accurately represent the whole counter array, no matter what kind of coding/decoding techniques are used.

### A. Worst-Case Average Binary Length of Counters

When no coding techniques are used, every counter value is stored as a plain binary number. The minimum memory required are the worst-case total number of bits of those counter values.



For each counter value  $C_i$ , its binary length is  $\lfloor \log_2 C_i \rfloor + 1$ , where we use the convention  $\log_2(0) = 0$ . Hence, the worst-case bound of the average<sup>4</sup> number of bits per counter, defined as  $B_0$ , is as follows:

$$\begin{aligned} B_0 &= \max_{\sum_{i=1}^N C_i \leq M} \sum_{i=1}^N \frac{\lfloor \log_2 C_i \rfloor + 1}{N} \\ &= \max_{\sum_{i=1}^N C_i = M} \sum_{i=1}^N \frac{\lfloor \log_2 C_i \rfloor + 1}{N}. \end{aligned}$$

It could be bounded through the Jensen's inequality

$$\begin{aligned} B_0 &\leq \sum_{i=1}^N \frac{\log_2 C_i + 1}{N} = \sum_{i=1}^N \frac{\log_2 C_i}{N} + 1 \\ &\leq \log_2 \left( \frac{1}{N} \sum_{i=1}^N C_i \right) + 1 = \log_2(M/N) + 1. \end{aligned}$$

If  $\log_2 \frac{M}{N}$  is an integer, all the inequalities above would hold when all counter values are equal to  $\frac{M}{N}$ , and the bound  $B_0$  would be reached at  $\log_2 \frac{M}{N} + 1$ .

When  $\log_2 \frac{M}{N}$  is not an integer, it becomes a little more complex to get the precise value of  $B_0$ . We claim  $B_0$  would be bounded by the result of the following optimization:

$$\begin{aligned} B_0 &= \max \sum_{i=1}^L i\beta_i \\ \text{subject to } &\forall i, \beta_i \geq 0; \sum_{i=1}^L \beta_i = 1 \end{aligned} \quad (7)$$

$$M/N + 1 \leq \sum_{i=1}^L \beta_i 2^i \leq 2M/N \quad (8)$$

where  $\beta_i$  is the fraction of  $i$ -b-long counters among all counters, i.e.,  $(1/N) \sum_{j=1}^N 1_{\lfloor \log_2 C_j \rfloor = i-1}$ ,  $L$  is the maximum possible length of one counter, i.e.,  $\lfloor \log_2 M \rfloor + 1$ . The constraints (8) are transformed from the constraint  $\sum_{i=1}^N C_i = M$ , since for each  $C_i$  that is  $j$  b long,  $2^{j-1} \leq C_i \leq 2^j - 1$ .

We could numerically solve the linear programming above and could find that the bound  $B_0$  is at least  $\log_2 \frac{M}{N} + 0.9$  for arbitrary  $\frac{M}{N}$ .

#### B. Bound<sub>1</sub>: Considering the Unavoidable Indexing Costs

The lower bound above does not account for the extra bits needed to delimit these counter values. However, it is hard to directly get the lower bound of indexing cost since we would never know whether we have invented the most efficient indexing scheme and how far it is from the optimal cost. We can only approximately estimate the cost through information entropy.

Given an array of indistinguishable counter bits, the original counter values could be decoded if and only if the sequence of the sizes of the counters are known. Hence, we could model the worst-case indexing cost by calculating the worst-case information entropy of the sequence of the counter sizes under all possible distributions as follows:

$$B_{\text{index}} = \max_{\text{subject to (7),(8)}} \left\{ -\sum_{i=1}^L \beta_i \log_2 \beta_i \right\}.$$

<sup>4</sup>For the convenience of comparison, we calculate and compare the *average* number of bits per counter instead of the total number.

TABLE I  
BOUNDS FOR SCHEMES WITHOUT USING ANY CODING TECHNIQUE

$\log_2(\frac{M}{N}) =$	2	4	6	8
$B_0 = \log_2 \frac{M}{N} +$	1.00	1.00	1.00	1.00
$B_{\text{index}} =$	1.83	2.61	3.10	3.45
$B_1 = \log_2 \frac{M}{N} +$	1.82	1.94	1.96	1.97

However, we should notice that the counter size distribution  $\{\beta_1, \beta_2, \dots\}$  for the worst-case indexing cost may not be the same as the distribution for the worst-case counter bits. Hence, if we want a more reasonable bound, it should be the result of the following optimization of the sum of both the indexing entropy and the counter bits:

$$B_1 = \max_{\text{subject to (7),(8)}} \left\{ -\sum_{i=1}^L \beta_i \log_2 \beta_i + \sum_{i=1}^L i\beta_i \right\}.$$

Both bounds could be solved by standard Lagrange techniques. The numerical result of  $B_0$ ,  $B_{\text{index}}$  and  $B_1$  are presented in Table I. Although the worst-case indexing cost  $B_{\text{index}}$  could be very large, only bound  $B_0$  and  $B_1$  are comparable to our scheme. Compared to numbers in Table III, our schemes are about 3–4 b from the optimal cost that one could achieve without compressing the original counter bits.

#### C. Bound<sub>2</sub>: Lower Bound When Optimal Coding is Used

Although we have not seen any works that could use coding techniques and support router-level fast random read and write at the same time, we are still interested in the optimal memory cost if we allow coding techniques. In the conference version of this paper [11], we calculate the worst-case 2-D empirical information entropy of the whole counter array and claim that it is the worst-case bound for the case when optimal coding is used. However, the method employed by [11] suffers the similar weakness as  $B_1$  and hence is not strict.

We propose the following method to calculate a strict worst-case bound. Our constraint  $\sum_{i=1}^N C_i \leq M$  is equivalent to  $\sum_{i=1}^N C_i + \delta = M$ , where  $\delta$  is a nonnegative integer. Basic combinatorics gives that the number of distinct nonnegative integer vectors  $\{C_1, \dots, C_N, \delta\}$  satisfying  $\sum_{i=1}^N C_i + \delta = M$  is exactly  $\binom{M+N}{N}$ . Since each possible counter value vector must correspond to a distinct memory state, at least  $\log_2 \binom{M+N}{N}$  b of memory are needed to accurately represent all possible counter value vectors, no matter what coding technique is employed. Therefore, we get the lower bound on average number of bits per counter as  $B_2 = \log_2 \binom{M+N}{N} / N$ . This lower bound is also achievable in theory since we can index all possible counter value vectors from 1 through  $\binom{M+N}{N}$  and simply store the index as a binary number.

Using Sterling's Formula, we could get

$$\begin{aligned} B_2 &= \log_2 \binom{M+N}{N} / N \approx \log_2 \left( \frac{M}{N} \left( 1 + \frac{N}{M} \right)^{\frac{M+N}{N}} \right) \\ &\approx \left( \log_2 \frac{M}{N} + 1.44 \right). \end{aligned}$$

The numerical values of  $B_2$  are presented in Table II. Interestingly, the results are very close to the worst-case empirical entropy results in [11], with differences less than  $1 \times 10^{-4}$ .

TABLE II  
INFORMATION-THEORETIC LOWER BOUND

$\log_2(\frac{M}{N}) =$	2	4	6	8
$B_2 = \log_2 \frac{M}{N} +$	1.61	1.49	1.45	1.45

We observe that the worst-case bound  $B_2$ , which allows coding techniques to be used on counter values, would be smaller than  $B_1$  by around 0.2–0.5 b. Hence, we could conclude that we are only 4–6 b away from the optimal cost we could get even if optimal coding is used.

We should notice that in reality we will never be even close to the bounds above, even the bound  $B_1$ , because the information theory employed in all those calculations does not take account of the complexity issues, considering our application scenarios require encoding and decoding (write and read) to be both very fast.

## V. PERFORMANCE EVALUATIONS

In this section, we will evaluate the performance of BRICK and show that BRICK is extremely memory-efficient. Our results show that the number of extra bits needed per counter in addition to the lower bounds  $\log_2 \frac{M}{N}$  remains practically constant with increasing number of flows  $N$ , and hence the solution is scalable. We also evaluate in Section V-B the performance of our architecture using two real-world Internet traffic traces. Finally, we discuss implementation issues in Section VI.

### A. Numerical Results of Analytical Bounds

In this section, we present a set of numerical results computed from the tail bound theorems derived in Section III.

1) *Configuration of Parameters and Memory Costs Optimization*: Recall that in our problem, the number of flows  $N$  and the maximum total increments in a measurement period  $M$  are given. For the specified  $N$  and  $M$ , we apply our tail bound theorems to derive optimal configurations for different combinations of constraints, i.e., bucket sizes  $k$ , number of levels  $p$ , and failure probabilities  $P_f$ . The derived configuration is in terms of the number of entries in each subcounter array  $k_j$ , the width of each subcounter array  $w_j$ , and the number of full-size buckets  $J$  that we need to ensure a failure probability less than  $P_f$  (the probability that we have insufficient entries in a subcounter array or a full-size bucket).

For a configuration of these parameters, the amount of memory required, which is also the target of our optimization, can be computed as follows:

$$\mathcal{S}_\ell = \left( \left\lceil \sum_{j=1}^p k_j(w_j + 1) \right\rceil - k_p \right) + (\lfloor \log_2 J \rfloor + 2) \quad (9)$$

$$\mathcal{S} = \left( \sum_{\ell=1}^h \mathcal{S}_\ell \right) + Jk(L + 1). \quad (10)$$

Here,  $\mathcal{S}_\ell$  is the memory cost of each bucket; its first component corresponds to the space required for storing the subcounter arrays and the index bitmaps, and its second component corresponds to the overflow status flag and the index to the corre-

TABLE III  
EXAMPLE OF SUBCOUNTER ARRAY SIZING AND PER-COUNTER STORAGE FOR  $k = 64$  AND  $P_f = 10^{-10}$ . (a) SIZING OF SUBCOUNTER ARRAYS. (b) SIZE OF EACH SUBCOUNTER ARRAY =  $k_j \times w_j$  (IN BITS). (c) STORAGE PER COUNTER

$p$	$k_2$	$k_3$	$k_4$	$k_5$	$w_1$	$w_2$	$w_3$	$w_4$	$w_5$
3	15	3			$\log_2 \frac{M}{N} + 3$	4	13		
4	25	10	2		$\log_2 \frac{M}{N} + 2$	2	4	12	
5	25	10	3	1	$\log_2 \frac{M}{N} + 2$	2	3	4	9

(a)

$p$	$A_2$	$A_3$	$A_4$	$A_5$
3	$15 \times 4 = 60$	$3 \times 13 = 39$		
4	$25 \times 2 = 50$	$10 \times 4 = 40$	$2 \times 12 = 24$	
5	$25 \times 2 = 50$	$10 \times 3 = 30$	$3 \times 4 = 12$	$1 \times 9 = 9$

(b)

$p = 3$	$p = 4$	$p = 5$
$\log_2 \frac{M}{N} + 6.05$	$\log_2 \frac{M}{N} + 5.66$	$\log_2 \frac{M}{N} + 5.50$

(c)

sponding full-size bucket.<sup>5</sup> Then, the total memory cost  $\mathcal{S}$  is  $h = \lceil N/k \rceil$  buckets of size  $\mathcal{S}_\ell$  each plus  $J$  full-size buckets of size  $k(L + 1)$  each. (For each full-size counter of size  $L$ , we need 1 b for indicating the migration status.)

We traverse typical possible configurations with the help of branch-and-prune technique and select the best among those configurations. We should emphasize that there might actually exist many different configurations around the optimal point. The difference between the average memory utilization of those quasi-optimal points is very minute. In practice, it is not necessary to really find the “optimal” point as long as the solution found is close to the optimal enough. Hence, the results shown are only of typical examples and might not be the best.

2) *Numerical Results With Various Configurations*: In Section IV, we have noticed that all lower bounds are just around 1 or 2 b over  $\log_2 \frac{M}{N}$ . Hence, for convenience, we use  $\sigma = (\mathcal{S}/N) - \log_2(M/N)$  as a metric to evaluate the space efficiency of our solution.

In Table III, we first consider results for the case with  $k = 64$  counters per bucket. We first consider this case with a small bucket size to ensure that all string operations are within 64 b, which allows for direct implementations using 64-b instructions in modern processors [1], [2]. As we shall see, substantial statistical multiplexing can already be achieved with  $k = 64$ . For the results presented in Table III, we used representative parameters with  $N = 1$  million counters and  $M = 16$  million as the maximum total increments during a measurement period. We also set the failure probability to be  $P_f = 10^{-10}$ , which is a tiny probability corresponding to an average of one failure (when there are more than  $J$  overflowed buckets) every 10 000 years. We will later show in Figs. 5–7 that the additional per-counter storage cost beyond the minimum width of the average count is practically a constant unrelated to the number of flows  $N$ , the maximum total increments  $M$ , or the failure probability  $P_f$ .

In Table III(a), the number of entries and the width for each subcounter array are shown for BRICK implementations with varying number of levels  $p$ . As can be seen, in each design, the number of entries decreases exponentially as we go to the higher subcounter arrays. This is the main source of our compression.

<sup>5</sup>Since there are  $J$  full-size buckets, this index can be stored in  $\lfloor \log_2 J \rfloor + 1$  b.

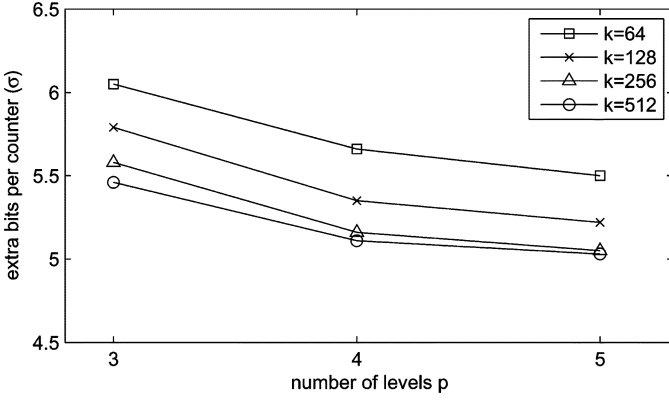


Fig. 4. Impact of increasing bucket size  $k$ . Extra bits  $\sigma$  in the range of [5.03, 6.05].

With  $k = 64$ , the rank indexing operation described in Section II only needs to be performed on bitmaps with  $|I_j| \leq 64$  b (much less than 64 for the higher subcounter arrays) and can be directly implemented using 64-b popcount and bitwise-logical instructions that are available in modern processors [1], [2]. Table III(b) shows the size of each subcounter array. For all three designs, the space requirement for each subcounter array other than  $A_1$  is also less than 64 b. Therefore, the “varshift” operator described in Section II-C, which only needs to operate on  $A_2$  and higher, can be directly implemented using 64-b shift and bitwise-logical instructions as well.

In Table III(c), the per-counter storage cost for the three designs are shown. For three levels, the extra storage cost per counter is 6.05, and the extra storage costs per counter are 5.66 and 5.50 for four and five levels, respectively. The amount of extra storage only decreases slightly with additional levels in the BRICK implementation. For example, as we go from three to five levels, the reduction of  $6.05 - 5.50 = 0.55$  extra bits is only about 5.5% in the overall per-counter cost if  $\log_2(M/N) = 4$ .

We next consider the impact of larger bucket sizes on storage costs. Fig. 4 shows the results for  $k = 128, 256$ , and  $512$ . The number of extra bits per counter decreases with increasing bucket sizes and number of levels, with  $\sigma$  in the range of [5.03, 6.05]. The results show that increasing the bucket size has only an insignificant impact on the storage savings, corresponding to only a small increase in statistical multiplexing with larger buckets. Therefore, we will use 64 counters per bucket for software implementation and recommend it for ASIC implementation as well since it has the advantage that operations can be directly implemented using 64-b processor instructions in software.

As stated earlier, the added per-counter cost is practically a constant with respect to the number of flows  $N$ , the maximum total increments  $M$ , and the failure probability  $P_f$ . We now verify this statement by evaluating BRICK under different  $N$ ,  $M$ , and  $P_f$ . The results are shown in Figs. 5–7, respectively. In these evaluations, we used 64 counters per bucket and four levels.

In Figs. 5–7, we evaluate the impact of different  $N$ ,  $M$ , and  $P_f$ , where we use  $k = 64$  and  $p = 4$ .

Fig. 5 shows that the added per-counter cost remains practically constant as we increase  $N$  exponentially by powers of 10.

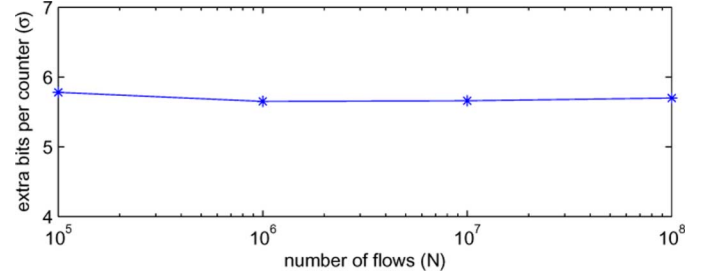


Fig. 5. Impact of increasing number of flows  $N$ . Extra bits  $\sigma$  in the range of [5.65, 5.78].

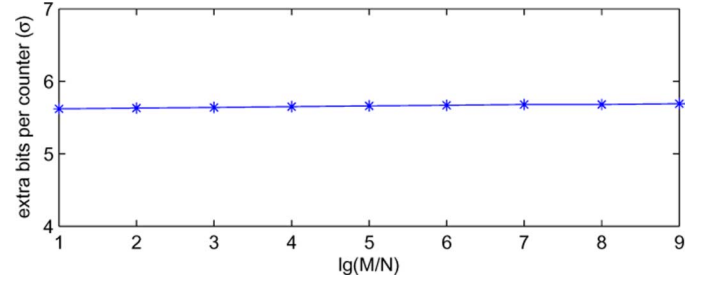


Fig. 6. Impact of increasing  $\log_2(M/N)$ . Extra bits  $\sigma$  in the range of [5.62, 5.69].

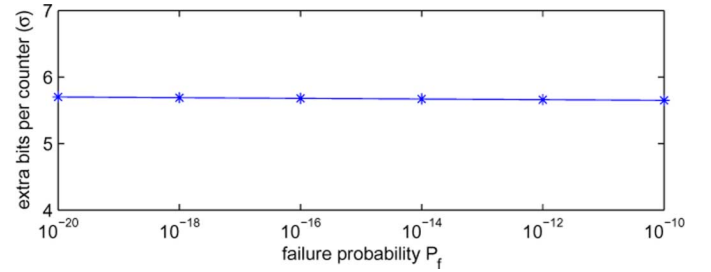


Fig. 7. Impact of decreasing failure probability  $P_f$ . Extra bits  $\sigma$  in the range of [5.65, 5.70].

Similarly, Fig. 6 shows that the added cost also remains practically constant with different ratios of  $M$  and  $N$ . These results show that BRICK is scalable to different values of  $M$  and  $N$  with per-counter storage cost within approximately a constant factor from the minimum width of the average count. Fig. 7 shows the impact of decreasing failure probability. We show results for  $P_f = 10^{-10}$  down to  $10^{-20}$ . Again, we see that the change in storage cost is negligible with decreasing failure probability, which means BRICK can be optimized to vanishingly small failure probabilities with virtually no impact on storage cost.

### B. Results for Real Internet Traces

In this section, we evaluate our active counter architecture using real-world Internet traffic traces. The traces that we used were collected at different locations in the Internet, namely the University of Southern California (USC) and University of North Carolina (UNC), respectively. The trace from USC was collected at their Los Nettos tracing facility on February 2, 2004, and the trace from UNC was collected on a 1-Gb/s access link connecting the campus to the rest of the Internet on April 24, 2003. For each trace, we used a 10-min segment,

TABLE IV  
PERCENTAGE OF FULL-SIZE BUCKETS

Trace	$h$	$J$	$\frac{J}{h}$	$J^*$	$J^{**}$
USC	17.3K	112	0.65%	99	0
UNC	19.5K	127	0.65%	172	0

corresponding to a measurement epoch. The trace segment from USC has 18.9 million packets and around 1.1 million flows; the trace segment from UNC has 32.6 million packets and around 1.24 million flows. We use the first counter for the first encountered flow in the trace, the second counter for the second encountered flow in the trace, etc.

We use the same general parameter settings as the evaluations in Section IV with 64 counters per bucket, four levels, and a failure probability of  $P_f = 10^{-10}$ . The total storage space required for counting packets<sup>6</sup> in USC trace is 1.36 MB, and the total required for the UNC trace is 1.61 MB. In comparison, a naive implementation would require a worst-case counter width for all counters. Both traces require a worst-case width of 25 b, whereas the BRICK implementations require a per-counter cost of about 10 b. The total storage required for a naive implementation is 3.85 MB for the USC trace and 4.40 MB for the UNC trace. The BRICK implementations represent a 2.5 $\times$  improvement in both cases. This is exciting since, with the same amount of memory, we will be able to squeeze in 2.5 times more counters, which is badly needed in future faster and “more crowded” Internet.

Table IV shows the number of full-size buckets needed according to our tail bounds, and the number of full-size buckets actually used. We should emphasize that the numbers shown in Table IV are only one example of the various good configurations selected by the process described in Section V-A1. The total number of full-size buckets needed could be much smaller by two to 10 folds, at the cost of a little more total memory utilization (typically around 1 to 5%).

In Table IV, we see that only a small number of full-size buckets are needed to guarantee a tiny probability ( $P_f = 10^{-10}$ ) that we will have insufficient number of full-size buckets to handle bucket overflows.  $J^*$  denotes the actual number of full-size buckets used when no random permutation is used.  $J^{**}$  denotes the actual number of full-size buckets used when permutation is done by reversing bits. We could see that  $J^*$  is similar to or even worse than the calculated  $J$ . This is because our guarantee is based on randomized permutation, and no permutation is used for  $J^*$ . In this experiment, the larger flows tend to be concentrated in lower index counters, thus causing those buckets to overflow to full-size buckets. However, when some very simple “randomization” techniques are used, such as reversing the bits of the counter index, we could see that actually no full-size buckets are used for both traces. We emphasize that the  $J$  full-size buckets are allocated for guaranteeing a tiny probability ( $P_f = 10^{-10}$ ) of overflow for any counter value distribution. In summary, this experiment demonstrates random permutations are crucial to our design and to establishing our statistical guarantee and verifies that simple permutation techniques might be sufficient for real-world deployment.

<sup>6</sup>We note that the memory savings for counting bytes would be less, due to the much larger  $M$ . Considering the typical average packet size is around 500 B, 9 more bits are needed per counter for both the naive implementation and the BRICK implementation.

## VI. DISCUSSIONS

### A. Implementation Issues

In a BRICK implementation, all subcounter arrays ( $A_j$ ) and index bitmaps ( $I_j$ ) are fixed in size, and the number and size of buckets are also fixed. Consider the three-level case shown in Table III with  $k = 64$ . Both lookup and increment operations can be performed with 10 memory accesses in total, five reads and five writes. For the bucket being read or updated, we first retrieve all bitmaps ( $I_j$ ), bucket overflow status flag  $f_\ell$ , and an index field  $t_\ell$  to a full-size bucket in case a bucket overflow has previously occurred. All this information for a bucket can be retrieved in two memory reads with 64-b words, the first word corresponds to  $I_1$  with 64 b, and the second word stores  $I_2 = 3$  b, the overflow status flag, and the  $t_\ell$  (about 7 b). If  $f_\ell$  is not set, then we need up to three reads and writes to update the three levels of subcounter arrays. The updated index bitmaps and overflow status flags can be written back in two memory writes. If  $f_\ell$  has been set, then we read directly from the corresponding entry in the full-size bucket indicated by  $t_\ell$  for a lookup operation, avoiding the need to read the subcounter arrays, hence requiring fewer memory accesses. Similarly, an increment operation for a counter that is already in a full-size bucket takes only one read and one write to update. If a bucket overflow occurs during an increment of a counter in a bucket, there is no need to access the last subcounter array (otherwise, we would not have an overflow). Therefore, we save two memory accesses at the expense of one write to the full-size bucket. With index bitmaps, overflow status flag, and full-size index field packed into two words, the worst-case number of memory accesses is 10 in total, which permits updates in 20 ns with a 2-ns SRAM time, enabling over 15 million packets per second of updates.

BRICK is also amenable to pipelining in hardware. In general, the lookup or update of each successive level of subcounter arrays can be pipelined such that at each packet arrival, a lookup or update can operate on  $A_1$  while a previous operation operates on  $A_2$ , and so forth. This enables the processing of hundreds of millions of packets per second.

## VII. CONCLUSION

We presented a novel exact active statistics counter architecture called BRICK (Bucketized Rank Indexed Counters) that can very efficiently store large arrays of variable-width counters entirely in SRAM while supporting extremely fast increments and lookups. This high memory (SRAM) efficiency is achieved through a statistical multiplexing technique, which by grouping a fixed number of randomly selected counters into a bucket, allows us to tightly bound the amount of memory that needs to be allocated to each bucket. Statistical guarantees of BRICK are proven using a combination of stochastic ordering theory and probabilistic tail bound techniques. We also developed an extremely simple and memory-efficient indexing structure called rank indexing to allow for fast random access of every counter inside a bucket. Experiments with real-world Internet traffic traces show that our solution can indeed maintain large arrays of exact active statistics counters with moderate amounts of SRAM.

## REFERENCES

- [1] "Intel 64 and IA-32 architectures software developer's manual," Intel, Santa Clara, CA, 2007, vol. 2B [Online]. Available: <http://download.intel.com/technology/architecture/new-instructions-paper.pdf>
- [2] "Software optimization guide for AMD family 10h processors AMD, Sunnyvale, CA, 2007 [Online]. Available: [http://www.amd.com/us-en/assets/content\\_type/white\\_papers\\_and\\_tech\\_docs/40546.pdf](http://www.amd.com/us-en/assets/content_type/white_papers_and_tech_docs/40546.pdf)
- [3] S. Cohen and Y. Matias, "Spectral bloom filters," in *Proc. ACM SIGMOD*, 2003, pp. 241–252.
- [4] G. Cormode and S. Muthukrishnan, "An improved data stream summary: The count-min sketch and its applications," *J. Algor.*, 2004.
- [5] T. M. Cover and J. A. Thomas, *Information Theory*, 2nd ed. Hoboken, NJ: Wiley, 2005, ch. 5.6.
- [6] A. Cvetkovski, "An algorithm for approximate counting using limited memory resources," in *Proc. ACM SIGMETRICS*, 2007, pp. 181–190.
- [7] M. Degermark, A. Brodnik, S. Carlsson, and S. Pink, "Small forwarding tables for fast routing lookups," in *Proc. ACM SIGCOMM*, 1997, pp. 3–14.
- [8] C. Estan and G. Varghese, "New directions in traffic measurement and accounting," in *Proc. ACM SIGCOMM*, Aug. 2002, pp. 323–336.
- [9] R. G. Gallager, *Low-Density Parity-Check Codes*. Cambridge, MA: MIT Press, 1963.
- [10] R. GöbDoi, "Bounds for median and 50 percentage point of binomial and negative binomial distribution," *Metrika*, vol. 41, no. 1, pp. 43–54, 2003.
- [11] N. Hua, B. Lin, J. J. Xu, and H. C. Zhao, "Brick: A novel exact active statistics counter architecture," in *Proc. ANCS*, 2008, pp. 89–98.
- [12] N. Hua, H. Zhao, B. Lin, and J. Xu, "Rank-indexed hashing: A compact construction of bloom filters and variants," in *Proc. IEEE ICNP*, Oct. 2008, pp. 73–82.
- [13] G. Jacobson, "Space-efficient static trees and graphs," in *Proc. 30th FOCS*, 1989, pp. 549–554.
- [14] B. Krishnamurthy, S. Sen, Y. Zhang, and Y. Chen, "Sketch-based change detection: Methods, evaluation, and applications," in *Proc. ACM SIGCOMM IMC*, Oct. 2003, pp. 234–247.
- [15] A. Kumar, M. Sung, J. Xu, and J. Wang, "Data streaming algorithms for efficient and accurate estimation of flow size distribution," in *Proc. ACM SIGMETRICS*, 2004, pp. 177–188.
- [16] Y. Lu, A. Montanari, B. Prabhakar, S. Dharmapurikar, and A. Kabbani, "Counter braids: A novel counter architecture for per-flow measurement," in *Proc. ACM SIGMETRICS*, 2008, pp. 121–132.
- [17] M. Mitzenmacher and E. Upfal, *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge, U.K.: Cambridge Univ. Press, 2005.
- [18] R. Morris, "Counting large numbers of events in small registers," *Commun. ACM*, vol. 21, no. 10, pp. 840–842, Oct. 1978.
- [19] R. Motwani and P. Raghavan, *Randomized Algorithms*. Cambridge, U.K.: Cambridge Univ. Press, 1995.
- [20] A. Müller and D. Stoyan, *Comparison Methods for Stochastic Models and Risks*. Hoboken, NJ: Wiley, 2002.
- [21] S. Ramabhadran and G. Varghese, "Efficient implementation of a statistics counter architecture," in *Proc. ACM SIGMETRICS*, Jun. 2003, pp. 261–271.
- [22] M. Roeder and B. Lin, "Maintaining exact statistics counters with a multilevel counter memory," in *Proc. IEEE GLOBECOM*, Dallas, TX, 2004, vol. 2, pp. 576–581.
- [23] D. Shah, S. Iyer, B. Prabhakar, and N. McKeown, "Maintaining statistics counters in router line cards," *IEEE Micro*, vol. 22, no. 1, pp. 76–81, Jan.–Feb. 2002.
- [24] R. Stanojevic, "Small active counters," in *Proc. IEEE INFOCOM*, 2007, pp. 2153–2161.
- [25] N. Tuck, T. Sherwood, B. Calder, and G. Varghese, "Deterministic memory-efficient string matching algorithms for intrusion detection," in *Proc. IEEE INFOCOM*, 2004, vol. 4, pp. 2628–2639.
- [26] Y. Zhang, S. Singh, S. Sen, N. Duffield, and C. Lund, "Online identification of hierarchical heavy hitters: Algorithms, evaluation, and application," in *Proc. ACM SIGCOMM IMC*, Oct. 2004, pp. 101–114.
- [27] Q. Zhao, A. Kumar, J. Wang, and J. Xu, "Data streaming algorithms for accurate and efficient measurement of traffic and flow matrices," in *Proc. ACM SIGMETRICS*, Jun. 2005, pp. 350–361.
- [28] Q. Zhao, J. Xu, and Z. Liu, "Design of a novel statistics counter architecture with optimal space and time efficiency," in *Proc. ACM SIGMETRICS*, Jun. 2006, pp. 323–334.



**Nan Hua** received the B.E. degree in information and electronic engineering and M.E. degree in information and telecommunication engineering from Tsinghua University, Beijing, China, in 2005 and 2007, respectively, and is currently a Ph.D. student with the College of Computing, Georgia Institute of Technology, Atlanta.

His current research interests include hardware algorithms and data structures for network routers and network measurements.



**Jun (Jim) Xu** received the Ph.D. degree in computer and information science from The Ohio State University, Columbus, in 2000.

He is an Associate Professor with the College of Computing, Georgia Institute of Technology, Atlanta. His current research interests include data-streaming algorithms for the measurement and monitoring of computer networks, hardware algorithms and data structures for network routers, network security, and performance modeling and simulation.

Dr. Xu received the NSF CAREER Award in 2003, the ACM SIGMETRICS Best Student Paper Award in 2004, and IBM faculty awards in 2006 and 2008.

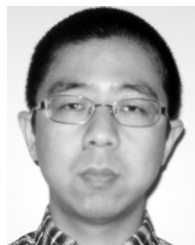


**Bill Lin** received the B.S., M.S., and Ph.D. degrees in electrical engineering and computer sciences from the University of California, Berkeley, in 1985, 1988, and 1991, respectively.

He is a Professor of electrical and computer engineering with the University of California, San Diego, where he is actively involved with the Center for Wireless Communications (CWC), the Center for Networked Systems (CNS), and the California Institute for Telecommunications and Information Technology (Calit2) in industry-sponsored research

efforts. Prior to joining the faculty at UCSD, he was the head of the System Control and Communications Group, IMEC, Leuven, Belgium. IMEC is the largest independent microelectronics and information technology research center in Europe. It is funded by European funding agencies in joint projects with major European telecom and semiconductor companies. His research has led to over 150 journal and conference publications. He also holds three awarded patents.

Prof. Lin has served on panels and given invited presentations at several major conferences. He has served on over 35 program committees, including serving as the General Chair for NOCS 2009, ANSC 2010, and IWQoS 2011. He has received two best paper awards, two best paper nominations, and two distinguished paper citations.



**Haiquan (Chuck) Zhao** received the Ph.D. degree in computer science from the Georgia Institute of Technology, Atlanta, in 2010.

He joined Microsoft, Redmond, WA, in 2010.