

HeavyKeeper: An Accurate Algorithm for Finding Top- k Elephant Flows

Tong Yang, Haowei Zhang, Jinyang Li, Junzhi Gong, Steve Uhlig, Shigang Chen and Xiaoming Li

Abstract—Finding top- k elephant flows is a critical task in network traffic measurement, with many applications in congestion control, anomaly detection and traffic engineering. As the line rates keep increasing in today’s networks, designing accurate and fast algorithms for online identification of elephant flows becomes more and more challenging. The prior algorithms are seriously limited in achieving accuracy under the constraints of heavy traffic and small on-chip memory in use. We observe that the basic strategies adopted by these algorithms either require significant space overhead to measure the sizes of all flows or incur significant inaccuracy when deciding which flows to keep track of. In this paper, we adopt a new strategy, called *count-with-exponential-decay*, to achieve space-accuracy balance by actively removing small flows through decaying, while minimizing the impact on large flows, so as to achieve high precision in finding top- k elephant flows. Moreover, the proposed algorithm called HeavyKeeper incurs small, constant processing overhead per packet and thus supports high line rates. Experimental results show that HeavyKeeper algorithm achieves 99.99% precision with a small memory size, and reduces the error by around 3 orders of magnitude on average compared to the state-of-the-art.

Index Terms—HeavyKeeper, Top- k , Sketch, Network measurements, Elephant flow

I. INTRODUCTION

A. Background and Motivation

Finding the largest k flows, also referred to as the top- k elephant flows, is a fundamental network management function, where a flow’s ID is usually defined as a combination of certain packet header fields, such as source IP address, destination IP address, source port, destination port, and protocol type, and the size of a flow is defined as the number of packets of the flow. Elephant flows contribute a large portion of network traffic. Many management applications can benefit from a function that can find them efficiently, such as congestion control by dynamically scheduling elephant flows [2], network capacity planning [3], anomaly detection [4], and caching of forwarding table entries [5]. Such a function not only is important in networking measurements [6]–[15], but also has applications beyond networking in areas such as data mining [16]–[18], information retrieval [19], databases [20], and security [21].

The preliminary version of this paper has been published in USENIX ATC 2018 [1].

Tong Yang (yangtongemail@gmail.com), Haowei Zhang, Jinyang Li, Junzhi Gong and Xiaoming Li are with Department of Computer and Science, Peking University, Beijing, China. Tong Yang is also with Shenzhen Graduate School, Peking University.

Steve Uhlig is with Queen Mary University of London, UK. Email: steve@eccs.qmul.ac.uk.

Shigang Chen is with the Department of Computer & Information Science Engineering, University of Florida, US. Email: sgchen@ufl.edu.

In real network traffic, it is well known that the distribution of flow sizes (the number of packets in a flow), is highly skewed [22]–[29], *i.e.*, the majority are mouse flows, while the minority are elephant flows. Most flows are small while a few flows are very large. The small flows are usually called *mouse flows*, while the large ones are called *elephant flows*.

Finding the top- k elephant flows (or top- k flows for short) in high-speed networks is a challenging task [30]. Extremely high line rates of modern networks make it practically impossible to accurately track the information of all flows. Consequently, approximate methods have been proposed in the literature and gained wide acceptance [24], [31]–[37]. In order to keep up with the line rates, these algorithms are expected to use on-chip memory such as SRAM whose latency is around 1ns [38], [39], in contrast to a latency of around 50ns when off-chip DRAM is used [39]. However, on-chip memory is small. Adding to the challenge, it is highly desirable to keep per-packet processing overhead small and constant, which helps pipelining.

Traditional solutions to finding the top- k flows follow two basic strategies: *count-all* and *admit-all-count-some*. The count-all strategy relies on a sketch (*e.g.*, CM sketch [24]) to measure the sizes of all flows, while using a min-heap to keep track of the top- k flows. For each incoming packet, it records the packet in the sketch and retrieves from the sketch an estimate \hat{n}_i for the size of the flow f_i that the packet belongs to. If \hat{n}_i is larger than the smallest flow size in the min-heap, it replaces the smallest flow in the heap by flow f_i . As a large sketch is needed to count all flows, these solutions are not memory efficient.

The *admit-all-count-some* strategy is adopted by Frequent [40], Lossy Counting [34], Space-Saving [32] and CSS [31]. These algorithms are similar to each other. To save memory, Space-Saving only maintains a data structure called Stream-Summary to count only some flows (*e.g.*, m flows). Each new flow will be inserted into the summary, replacing the smallest existing flow. The initial size of the new flow is set as $\hat{n}_{min} + 1$, where \hat{n}_{min} is the size of the smallest flow in the summary. By keeping m flows in the summary, the algorithm will report the largest k flows among them, where $m > k$. It assumes every new incoming flow is an elephant flow, and expels the smallest one in the summary to make room for the new one. But most flows are mouse flows. Such an assumption causes significant error, especially under tight memory (for a limited value of m).

In addition to the above two categories of algorithms for finding top- k flows, there are many recent works [41]–[44] introducing a lot of new strategies, and we divide them as the

third category. The Elastic sketch uses votes to decide whether a flow should be recorded or evicted; HeavyGuardian uses the strategy of exponential decay to address five typical measurement tasks; Cold Filter uses a two-layer filter to prevent mouse flows from entering some data structures (e.g., Space-Saving, the CM sketch); and Counter Tree uses the strategy of two-dimensional counter sharing and derives mathematical formulas to estimate flow sizes.

B. Our Proposed Solution

In this paper, we propose a new algorithm, HeavyKeeper, which uses the similar strategy introduced from [42], called *count-with-exponential-decay*. It keeps all elephant flows while drastically reducing space wasted on mouse flows. HeavyGuardian can handle five different tasks, but not including top- k elephant flows detection, while the algorithm we proposed just focuses on finding top- k elephant flows. HeavyKeeper uses multiple arrays, and thus can scale well while HeavyGuardian cannot.

Unlike *count-all*, our strategy only keeps track of a small number of flows. Unlike *admit-all-count-some*, we do not automatically admit new flows into our data structure and the vast majority of mouse flows will be by-passed. For a small number of mouse flows that do enter our data structure, they will decay away to make room for true elephants. The decay is not uniform for the flows in our data structure. The design of exponential decay is biased against small flows, and it has a smaller impact on larger flows. This design works extremely well with real traffic traces under small memory.

II. PRELIMINARIES

A. Problem Statement

Simply speaking, finding top- k flows refers to finding the largest k flows. Let $\mathcal{P} = \mathbb{P}_1, \mathbb{P}_2, \dots, \mathbb{P}_N$ be a network stream with N packets. Each packet \mathbb{P}_l ($1 \leq l \leq N$) belongs to a flow f_i , where $f_i \in \mathcal{F} = \{f_1, f_2, \dots, f_M\}$ and \mathcal{F} is the set of flows. Let n_i be the real flow size of flow f_i in \mathcal{P} . We order all flows (f_1, f_2, \dots, f_M) so that $n_1 \geq n_2 \geq \dots \geq n_M$.

Given an integer k and a network stream \mathcal{P} , the output of top- k is a list of k flows from \mathcal{F} with the largest flow sizes, i.e., f_1, f_2, \dots, f_k .

B. Prior Art and Limitations

The count-all strategy: As mentioned above, the *count-all* strategy uses sketches (such as the CM sketch [24] or the Count sketch [33]) to record the sizes of all flows, and uses a min-heap to keep track of the top- k flows, including the flow IDs and their flow sizes. Take the CM sketch as an example. It records packets in a CM sketch, consisting of a pool of counters. For each arrival packet, it hashes the packet's flow ID f to d counters and increases these d counters by one. The smallest value of the d counters is used as the estimated size of the flow, which is used to update the min-heap.

The problem is that all flows are pseudo-randomly mapped to the same pool of counters through hashing. Each counter may be shared by multiple flows, and thus record the sum of sizes of all these flows. Consequently, a small flow may be

treated as an elephant flow if all its d counters are shared with real elephant flows.

The admit-all-count-some strategy: As mentioned above, quite a few algorithms use the *admit-all-count-some* strategy, including Frequent [40], Lossy counting [34], and Space-Saving [32]. Take Space-Saving as an example. It counts only the sizes of some flows in a data structure called Stream-Summary, which incurs $O(1)$ overhead to search a flow or update the smallest flow. For each arrival packet, if its flow ID is not in the summary, the flow will be admitted into the summary, replacing the smallest existing flow. The new flow's initial size is set to $\hat{n}_{min} + 1$, where \hat{n}_{min} is the smallest flow size in the summary before replacement. A recent work CSS [31] is proposed based on Space-Saving. It inherits the above strategy, but redesigns the data structure of Stream-Summary by using TinyTable [45] to reduce memory usage.

The strategy of *admit-all-count-some* is to admit all new flows while expelling the smallest existing ones from the summary. To give new flows a chance to stay in the summary, their initial flow sizes are set as $\hat{n}_{min} + 1$. Such a strategy drastically over-estimates sizes of flows, and we show an example here. Assume $\hat{n}_{min} = 10,000$ and the summary is already full. Given a new flow, it will directly replace the flow with the smallest size in the summary and set its size to be 10,001. If this new flow is a mouse flow, it is largely over-estimated. Therefore, numerous mouse flows will cause significant over-estimation errors.

III. THE DESIGN OF HEAVYKEEPER

In this section, we present the data structure and algorithm of our HeavyKeeper, and show how to find the top- k flows.

A. Rationale

We aim to use a small hash table to store all elephant flows. As there are a great number of flows, each bucket of the hash table will be mapped by many flows, and we aim to store only the largest flow with its size, which cannot be achieved with no error when using small memory. Therefore, we leverage a probabilistic method called *exponential-weakening decay*. Specifically, when the incoming flow is not found in the hashed bucket, we decay the flow size with a probability, which exponentially decreases as the flow size increases. If the flow size is decayed to 0, it replaces the original flow with the new flow. In this way, mouse flows can easily be decayed to 0, while elephant flows can easily keep stable in the bucket. There are two shortcomings: 1) With a small probability we elect the wrong flow as the largest flow; 2) The reported flow size might be under-estimated because of the decay operations. To address these problems, we use multiple hash tables with different hash functions. An elephant flow could be stored in multiple hash tables, we choose the recorded largest size, minimizing the error of flow sizes.

B. The HeavyKeeper Structure

As shown in Figure 1, HeavyKeeper is comprised of d arrays, and each array is comprised of w buckets. Each bucket

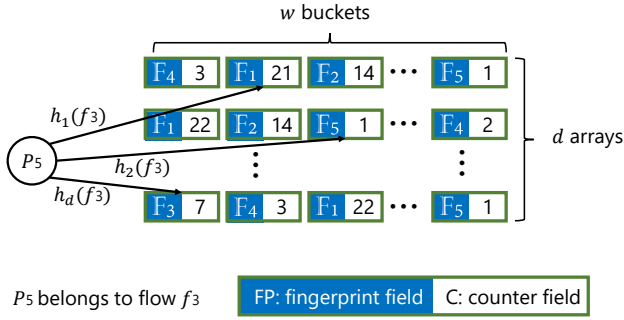


Fig. 1. The data structure of HeavyKeeper.

consists of two fields: a fingerprint field and a counter field.¹ For convenience, we use $A_j[t]$ to represent the t^{th} bucket in the j^{th} array, and use $A_j[t].FP$ and $A_j[t].C$ to represent its fingerprint field and counter field, respectively. Arrays $A_1 \dots A_d$ are associated with hash functions $h_1(\cdot) \dots h_d(\cdot)$, respectively. These d hash functions $h_1(\cdot) \dots h_d(\cdot)$ need to be 2-way independent.

Insertion: Initially, all fingerprint fields are *null*, and all counter fields are 0. For each incoming packet \mathbb{P}_l belonging to flow f_i , HeavyKeeper computes the d hash functions, and maps f_i to d buckets $A_j[h_j(f_i)]$ ($1 \leq j \leq d$) (one bucket in each array), which we call ***d mapped buckets*** for convenience. As shown in Figure 2, for each mapped bucket, HeavyKeeper applies different strategies for the following three cases:

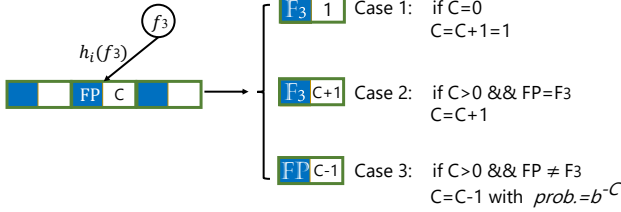


Fig. 2. The main insertion cases of HeavyKeeper. Note: 1) \mathbb{F}_3 is the fingerprint of flow f_3 . 2) $b > 1$ and $b \approx 1$ (e.g., $b = 1.08$). 3) In Case 3, when C is decayed to 0, the fingerprint field will be replaced by \mathbb{F}_3 , and then counter C is set to 1.

Case 1: When $A_j[h_j(f_i)].C = 0$. It means that no flow has been mapped to this bucket, then HeavyKeeper sets $A_j[h_j(f_i)].FP = \mathbb{F}_i$ and $A_j[h_j(f_i)].C = 1$, where \mathbb{F}_i represents the fingerprint of f_i .

Case 2: When $A_j[h_j(f_i)].C > 0$ and $A_j[h_j(f_i)].FP = \mathbb{F}_i$. It means $A_j[h_j(f_i)].C$ is possibly the estimated size of f_i . In this case, HeavyKeeper increments $A_j[h_j(f_i)].C$ by 1.

Case 3: When $A_j[h_j(f_i)].C > 0$ and $A_j[h_j(f_i)].FP \neq \mathbb{F}_i$. It means that $A_j[h_j(f_i)].C$ is not the estimated size of f_i . In here, HeavyKeeper applies the *exponential-weakening decay* strategy to this bucket: it decays $A_j[h_j(f_i)].C$ by 1 with a probability P_{decay} . After decay, if $A_j[h_j(f_i)].C = 0$, HeavyKeeper replaces $A_j[h_j(f_i)].FP$ with \mathbb{F}_i , and sets $A_j[h_j(f_i)].C$ to 1. Therefore, as long as flows are mapped to a bucket, its counter field will never be 0.

¹The fingerprint of a flow is a hash value generated by a certain function (for example, if we use $h_f(\cdot)$ as the fingerprint hash function, the fingerprint of flow f_j is $h_f(f_j)$). Although there can be hash collisions among flows, the probability is quite small. For example, if we set the fingerprint size to 16 bits, and there are 10000 buckets in the array, the probability of fingerprint collisions is $1.52 * 10^{-3}$.

Note that at any time the values of counters are non-negative, since decay only happens in **Case 3** and **Case 3** happens only when the value of the counter is larger than 0. And in **Case 3**, when a counter is decayed to zero, the new flow is inserted to this bucket and the counter is set to be 1 immediately.

Query: To query the size of a flow f_i , HeavyKeeper first computes the d hash functions to get d buckets $A_j[h_j(f_i)]$ ($1 \leq j \leq d$). Among the d mapped buckets, it chooses those buckets whose fingerprint fields are equal to \mathbb{F}_i . It then reports the maximum counter field of those buckets, i.e., $\max_{1 \leq j \leq d} \{A_j[h_j(f_i)].C\}$ where $A_j[h_j(f_i)].FP = \mathbb{F}_i$.

For convenience, for those d mapped buckets of f_i , if $A_j[h_j(f_i)].FP = \mathbb{F}_i$, we say that f_i is **held** at bucket $A_j[h_j(f_i)]$. Ignoring the limited impact of fingerprint collisions, we prove that *the reported size for each flow is equal to or smaller than the real flow size* in Section B. If a flow is held at no mapped bucket, it reports that it is a mouse flow. If a flow is held at multiple buckets, HeavyKeeper reports the maximum counter field.

Decay probability:

The key problem is how to choose a function to calculate the probability. Based on our experimental results on real and synthetic datasets, we find that as long as the parameters are set reasonably, functions satisfying the following condition all have a good performance: the larger the value in the current counter field is, the smaller the probability is. We finally choose the exponential function

$$P_{decay} = b^{-C} \quad (b > 1)$$

where C is the value in the current counter field and b ($b > 1$ and $b \approx 1$, e.g., $b = 1.08$) is a predefined exponential base number. This is because the function has the following properties. 1) As the value increases, the rate of probability reduction gradually increases and maps to $[0,1]$. 2) When the value is large enough (e.g., 50), the probability is close to 0, so we can regard the probability as 0, so as to accelerate the throughput of our algorithm. 3) When the value is small (e.g., 3), the recorded flow can hardly be an elephant flow, and at the same time the probability is close to 1, which exactly matches this condition.

Indeed, there are many other functions, which have a good performance, such as C^{-b} , $\frac{e^C}{1+e^C}$, etc. We have conducted experiments to compare those functions, and the experimental results show that the performances are similar with different decay functions.

Therefore, the larger size a flow has, the harder it is to decay its size. For elephant flows, it is held at several buckets, and the corresponding counter fields are incremented regularly, while decayed with a very small probability. Therefore, the error rate for estimated sizes of elephant flows is very small.

Note: Our data structure of d arrays and d 2-way independent hash functions may show some similarity with that of CM [24]. But similarity stops there. CM records the sizes of all flows; we record the sizes of a small number of flows. CM does not store flow IDs; we do. CM stores information of each flow in d counters; we keep each flow mostly in one bucket, while d -hashing helps find an empty bucket. CM does not have

to worry about the issue of kicking out existing flows to make room for new ones, which is what our exponential delay does.

Example: As shown in Figure 1, given an incoming packet \mathbb{P}_5 belonging to flow f_3 , we compute the d hash functions to obtain one bucket in each array. In the mapped bucket of the first array, the fingerprint field is not equal to \mathbb{F}_3 and the counter field is 21, thus we decay the counter field from 21 to 20 with a probability of 1.08^{-21} (assume $b = 1.08$). In the second mapped bucket, the fingerprint field is not \mathbb{F}_3 either, and with a probability of 1.08^{-1} , we decay the counter field from 1 to 0. If the counter field is decayed to 0, we set the fingerprint field to \mathbb{F}_3 , and set the counter field to 1. In the last mapped bucket, the fingerprint field is \mathbb{F}_3 , we increment the counter field from 7 to 8.

Analysis: HeavyKeeper uses fingerprint to identify and keep elephant flows. If a mouse flow with a small flow size is held at a bucket, it will be replaced by other flows mapped to this bucket soon, because each flow mapped to this bucket with a different fingerprint will decay the counter field with a high probability ($b^{-C} \rightarrow 1$ when C is small). If an elephant flow is held at a bucket, the corresponding counter field can easily be incremented to a large value since elephant flows have many incoming packets. Moreover, the decay probability becomes very small ($b^{-C} \rightarrow 0$ when C is large) as the counter field increases to a large value. Therefore, mouse flows can hardly be held in HeavyKeeper for a long time, and thus have a large probability to be *passers-by* of HeavyKeeper. However, elephant flows can keep stable in HeavyKeeper, and the estimated sizes of elephant flows are accurate.

C. Basic Version for Finding Top- k Elephant Flows

To find top- k elephant flows, our basic version just uses a HeavyKeeper and a min-heap. The min-heap is used to store the IDs and sizes of top- k flows. For each incoming packet \mathbb{P}_l belonging to flow f_i , we first insert it into HeavyKeeper. Suppose that HeavyKeeper reports the size of f_i as \hat{n}_i . If f_i is already in the min-heap, we update its estimated flow size with $\max(\hat{n}_i, \min_heap[f_i])$, where $\min_heap[f_i]$ is the recorded size of f_i in min-heap. Otherwise, if \hat{n}_i is larger than the smallest flow size which is in the root node of the min-heap, we expel the root node from the min-heap, and insert f_i with \hat{n}_i into the min-heap. To query top- k flows, we simply report the k flows in the min-heap with their estimated flow sizes.

Note that in our implementation, we use *Stream-Summary* instead of min-heap, as the function of min-heap and Stream-Summary is similar, and Stream-Summary can achieve $O(1)$ update complexity. For better understanding, we use min-heap to explain in our paper.

D. Optimizations

In this section, we propose further optimization methods to avoid accidental errors and improve speed. For convenience, we use n_{min} to denote the minimal flow size in the min-heap.

Optimization I: Fingerprint Collisions Detection.

Problem: Assume that there is a bucket in HeavyKeeper where flow f_i is held, and a mouse flow f_j mapped to the same bucket has the same fingerprint as f_i , i.e., $\mathbb{F}_i = \mathbb{F}_j$ due to hash collisions. Then, the mouse flow f_j is also held at this bucket,

and its estimated size is drastically over-estimated. In the worst case, if flow f_j has a fingerprint collision in all d arrays, the mouse flow f_j will probably be inserted into the min-heap. It can hardly be expelled due to its drastically over-estimated size. One effective solution is to store the entire IDs of flows instead of using fingerprints, which can definitely avoid hash collisions. However, in real data streams, the number of bits of a flow's ID is usually very large (e.g., more than 100 bits in 5-tuple), leading to a waste of memory. Indeed, the better the memory efficiency is, the higher the accuracy of algorithms will be. Our design goal is to find a solution to alleviate hash collisions without increasing the number of recorded bits. Therefore, our solution is to store fingerprints instead of entire IDs. In order to reduce the impact of hash collisions, we propose a solution based on the following Theorem.

Theorem 1. *When there is no fingerprint collision, after a flow f_i is inserted into HeavyKeeper, if its estimated size \hat{n}_i is larger than n_{min} (recall that we use n_{min} to denote the minimal flow size in the min-heap), then we must have*

$$\hat{n}_i = n_{min} + 1$$

The proof of this Theorem is not hard to derive and we skip it due to space limitations.

Solution: Based on Theorem 1, if f_i is not in the min-heap but $\hat{n}_i > n_{min} + 1$, then f_i is a mouse flow whose size is drastically over-estimated due to fingerprint collision. Therefore, we should not insert f_i into the min-heap in this case.

Optimization II: Selective Increment.

Problem: If a flow f_i is not in the min-heap, then the estimated flow size should be no larger than n_{min} . However, due to fingerprint collisions, there could be some mapped buckets of flow f_i where the fingerprint field is \mathbb{F}_i and the counter field is larger than n_{min} . In this case, flow f_i is not the flow that is held at this bucket, and thus increasing the corresponding counter field can only incur extra error.

Solution: In this case, instead of incrementing or decaying the corresponding counter field, we make no change.

E. Hardware Parallel Version

Based on the basic version, we propose a new version using the above two optimization methods. It is called *Hardware Parallel version* (Parallel version for short) because for each insertion, the operation in each array can be implemented in parallel on hardware platforms (e.g., FPGA, ASIC, or P4Switch). We will propose a more accurate version (named *Software Minimum version*, Minimum version for short in Section IV) at the cost of sacrificing the parallel property. The insertion and query processes of the Parallel version of our algorithm are presented as follows (see pseudo-code in Appendix A Algorithm 1 of our technical report [46]).

Insertion: All counters and fingerprints in HeavyKeeper and the min-heap are initialized to 0. For each incoming packet \mathbb{P}_l belonging to flow f_i , these are the following three steps for each insertion:

Step 1: We check whether flow f_i is already monitored by the min-heap, which is shown in line 1-3 in Appendix A

Algorithm 1. We use a boolean variable $flag$ to represent the result.

Step 2: We insert f_i into HeavyKeeper, which is shown in line 4-22 in Appendix A Algorithm 1. According to Optimization II, for each mapped bucket, if the fingerprint field is equal to \mathbb{F}_i , we increment the counter field only when $flag = true$ or $C < n_{min}$, where C is the original value in the counter field. *Step 3:* We get an estimated size \hat{n}_i of flow f_i from HeavyKeeper, which is shown in line 23-27 in Appendix A Algorithm 1. According to Optimization I, if $flag$ is *true*, we update the estimated size of flow f_i in the min-heap with \hat{n}_i . If $flag$ is *false*, we insert flow f_i into the min-heap with \hat{n}_i in only two cases: 1) the number of flows that are in the min-heap is less than k ; 2) $\hat{n}_i = n_{min} + 1$.

Query top- k flows: It reports the k flows recorded in the min-heap and their estimated flow sizes.

Analysis: Since HeavyKeeper achieves very small error rate on the flow size estimation of elephant flows, it can significantly reduce the error in finding top- k elephant flows. Furthermore, the first two optimizations reduce the impact of fingerprint collisions, and enhance the precision of finding top- k elephant flows and their flow size estimation.

F. Limitations and A Solution

As mentioned before, when the *exponential-weakening decay* is performed on a bucket, if its counter value is large enough (e.g., 50), the probability of reducing its value is close to 0. Therefore, in the worst case, when a new flow arrives, if all values of its mapped d counters are large enough, it could never be inserted into some buckets. In fact, this limitation means that the current memory size is too tight to record top- k elephant flows. To address this problem, we propose to use an extra global counter to record how many times this situation happens. As long as the value of the extra counter is larger than a predefined threshold, we add a new array, i.e., the $d + 1^{th}$ array. In this way, the new flow will have a chance to record its information.

Besides, our proposed algorithm cannot handle other flow measurement tasks (e.g., flow size estimation, entropy detection) and cannot support weighted updates. However, thanks to the fact that HeavyKeeper is designed mainly to handle top- k flows detection, it achieves higher accuracy than other related algorithms, which will be detailed in Section VI-E.

IV. SOFTWARE MINIMUM VERSION

In the above section, we describe the Hardware Parallel Version of HeavyKeeper, in which all the d arrays can be inserted or queried in parallel. We observe that its accuracy can be further improved by sacrificing the parallel property. In this section, we propose the Software Minimum Version to further enhance the accuracy.

A. Problem

We observe that it is unnecessary to decay all the mapped counters in the basic version. Specifically, when inserting an incoming packet \mathbb{P}_l belonging to flow f_i , HeavyKeeper computes d hash functions and maps f_i to d buckets $A_j[h_j(f_i)]$ ($1 \leq j \leq d$) (one bucket in each array). For each bucket,

HeavyKeeper applies different strategies for three different cases. We focus on the third case below. In **Case 3**, $A_j[h_j(f_i)].C > 0$ and $A_j[h_j(f_i)].FP \neq \mathbb{F}_i$, HeavyKeeper decays $A_j[h_j(f_i)].C$ by 1 with a probability P_{decay} , and after decay, if $A_j[h_j(f_i)].C = 0$, HeavyKeeper replaces $A_j[h_j(f_i)].FP$ with \mathbb{F}_i , and sets $A_j[h_j(f_i)].C$ to 1. However, for a bucket $A_k[h_k(f_i)]$ ($1 \leq k \leq d$) in HeavyKeeper where an elephant flow f_i is held, if another flow f_j is mapped to the same bucket due to hash collisions, i.e., $f_i \neq f_j$ and $\mathbb{F}_i = \mathbb{F}_j$, then $A_k[h_k(f_i)]$ is decayed by 1 with a probability P_{decay} , but such decay is not always necessary and could be harmful for the following reasons.

First, if f_j is a mouse flow which only has a few packets, the elephant flow f_i can hardly be replaced by it, but f_i 's counter field is possibly decayed for a few times (e.g., decayed from 1000 to 999). Such decay can hardly cause a replacement, but at the same time, it makes f_i 's recorded flow size in this bucket less than its real flow size, which will degrade the accuracy of queries.

Second, if f_j is an elephant flow which has a large number of packets, whether $A_k[h_k(f_i)].C$ will be decayed to 0 and $A_k[h_k(f_i)].FP$ will be replaced with \mathbb{F}_j depends on the following packets of f_i and f_j . In such a contest of the two elephant flows, the counter in this bucket may be decayed many times. There are two results. 1) If f_i wins and keeps held in this bucket, i.e., $A_k[h_k(f_i)].C$ never reaches 0, then $A_k[h_k(f_i)].C$ will be much less than the real flow size of f_i . When querying the size of flow f_i , HeavyKeeper reports the maximum counter field of all the mapped buckets. As an elephant flow, f_i is likely to be kept in several buckets, and the counter fields in other buckets may well be larger than $A_k[h_k(f_i)].C$, so $A_k[h_k(f_i)].C$ makes no contribution to the accuracy of queries. 2) If f_j wins and replaces f_i in the bucket $A_k[h_k(f_i)]$, after replacement, the counter starts from 1, so $A_k[h_k(f_i)].C$ is much less than the real flow size of f_j . Also, this counter makes no contribution to the query results of flow f_j . In summary, *it is unnecessary and unhelpful to decay large counters.*

It is possible that f_i will always occupy a bucket if we do not perform any decay on it. In the worst case, if f_i is not an elephant flow, this strategy will make new flows not have a choice to be inserted into that bucket. In other words, this method is not friendly to late-arrival elephant flows. However, this situation happens only when for a new flow, all values of its mapped d counters are very large. As mentioned in Section IV-C, we can use an extra counter and automatically add a new array to avoid this situation.

B. Solution: Minimum Decay

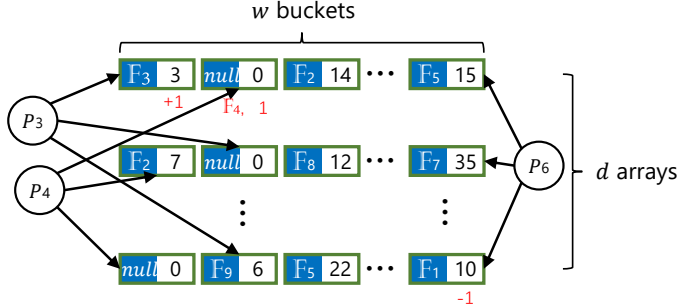
To address the above problem, we propose a solution, and the key technique is called ‘‘Minimum Decay’’. Its key idea is that we choose to decay the smallest one instead of decaying all the mapped counters. Below we show the details of our solution. For each incoming packet \mathbb{P}_l belonging to flow f_i , HeavyKeeper computes d hash functions and maps f_i to d buckets $A_j[h_j(f_i)]$ ($1 \leq j \leq d$) (one bucket in each array). For the d mapped buckets, suppose \mathbb{F}_i is the fingerprint of f_i . There are three situations.

Situation 1: If one of the d mapped buckets has the same fingerprint as \mathbb{F}_i , we just increment the corresponding counter by 1.

Situation 2: If all d mapped buckets do not have the fingerprint \mathbb{F}_i , but one or more of the mapped buckets are empty. In this situation, we just insert f_i into the first empty bucket.

Situation 3: If all d mapped buckets are full and do not have the fingerprint \mathbb{F}_i . In this situation, we choose the smallest counter among the mapped bucket, and then perform the decay operation. If there is more than one smallest counter, we only choose the first one to decay.

Note that for each insertion, we only update one mapped bucket, and do nothing for other mapped buckets.



P_3 belongs to flow f_3 : situation 1
 P_4 belongs to flow f_4 : situation 2
 P_6 belongs to flow f_6 : situation 3

FP: fingerprint field C: counter field

Fig. 3. Examples of the insertion of Parallel version.

Examples: Figure 3 shows three incoming packets corresponding to the three situations, respectively. Given each incoming packet, we compute the d hash functions to obtain one bucket in each array. We only show the first, second and last array for convenience. For packet \mathbb{P}_3 belonging to flow f_3 , the first mapped bucket holds the same fingerprint as f_3 (\mathbb{F}_3), so this is the above *Situation 1*. Thus we increment the counter field from 3 to 4. For packet \mathbb{P}_4 belonging to flow f_4 , none of the d mapped buckets holds the fingerprint \mathbb{F}_4 , but there are two empty buckets, so this is *Situation 2*. We insert flow f_4 into the mapped bucket in the first array. We set its fingerprint field to \mathbb{F}_4 and its counter field to 1. For packet \mathbb{P}_6 belonging to flow f_6 , none of the d mapped buckets holds the fingerprint \mathbb{F}_6 and none of them is full, so this is *Situation 3*. The counter field in the last mapped bucket is the smallest, so we decay it by 1 with a probability of 1.08^{-10} , and do nothing to the other mapped buckets.

C. Hardware Minimum Version for Finding Top-k Flow

Based on the Hardware Parallel version, we propose the Software Minimum version (Minimum version for short) using the above minimum decay technique. The insertion and query processes of our Minimum version of our algorithm are presented as follows. Due to space limitation, we present the pseudo-code in the Appendix of our technical report [46].

Insertion: All counters and fingerprints in HeavyKeeper and the min-heap are initialized to 0. For each incoming packet \mathbb{P}_i belonging to flow f_i , there are the following five steps for each insertion:

Step 1: We check whether flow f_i is already monitored by the min-heap, denoted by a boolean variable $flag$.

Step 2: We check whether there is a mapped bucket holding the same fingerprint as \mathbb{F}_i . If there is and the corresponding bucket could be updated ($flag = true$ or the value of counter is less than n_{min}), we increment the corresponding counter filed by 1, and then go to step 5; otherwise, we go to step 3.

Step 3: We check whether there is a mapped bucket that is empty. If there is, we insert this packet into the first empty bucket and then go to step 5; otherwise, we go to step 4.

Step 4: We choose the bucket with the smallest counter field among the d mapped buckets and decay it with a certain probability. If there is more than one such bucket, we only decay the first one.

Step 5: Step 5 is similar to step 3 of Parallel version of HeavyKeeper. We get an estimated size \hat{n}_i of flow f_i from HeavyKeeper. If $flag$ is *true*, we update the estimated size of flow f_i in the min-heap with \hat{n}_i . If $flag$ is *false*, we insert flow f_i into the min-heap with \hat{n}_i in only two cases: 1) the number of flows that are in the min-heap is less than k ; 2) $\hat{n}_i = n_{min} + 1$.

Query top-k flows: We report the k flows recorded in the min-heap and their estimated flow sizes.

Analysis: The Parallel version of HeavyKeeper achieves fast processing speed and small error rate in finding top-k elephant flows. Based on the Parallel version, the Minimum version further improves the accuracy. Specifically, when inserting a packet, the Minimum version only needs to change at most one bucket, thus it avoids unnecessary and unhelpful decay. Our experimental results (see Figure 23, 26 and 29) verify that the accuracy is significantly improved when using the Minimum Decay technique.

V. MATHEMATICAL ANALYSIS

In this section, we first claim that there is no over-estimation of HeavyKeeper, and then derive the formula of error bound in the Minimum version of HeavyKeeper. Note that we also derived the formula of error bound in the basic version of . Due to space limitation, we provide the derivation process of the basic version in the Appendix of our technical report [46].

A. Claim of No Over-estimation Error of HeavyKeeper

Theorem 2. In the Minimum version, let $n_i(t)$ be the real size of flow f_i at time t , and let $A_j[h_j(f_i)](t).C$ be the counter field of the mapped bucket of flow f_i in the j^{th} array at time t . If there is no fingerprint collision, then

$$\forall j, t, A_j[h_j(f_i)](t).C \leq n_i(t)$$

Proof. It is not hard to prove this theorem. Due to space limitation, we provide the proof in the Appendix of our technical report [46]. \square

B. Error Bound of the Minimum Version of HeavyKeeper

Theorem 3. Assume that there is no fingerprint collision and once the fingerprint of an elephant flow is inserted into its mapped bucket, it is held there all the time. For any $\epsilon > 0$, assume an elephant flow f_i with size n_i is held in the bucket, we have

$$Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} \leq \frac{\gamma}{\epsilon w n_i (b-1)} \quad (1)$$

where w is the width of each array, b the exponential base, and γ the proportion of mouse flows in all flows.

Proof. For convenience, we use N to denote the total number of packets, M to denote the number of different flows and d to denote the number of arrays. Let's focus on the j^{th} array. Flow f_i is correctly reported, so at the end, the fingerprint of flow f_i is held in the $h_j(f_i)^{\text{th}}$ bucket of the j^{th} array. Let $I_{i,j,i'}$ be a binary random variable, defined as

$$I_{i,j,i'} = \begin{cases} 0 & (f_i = f_{i'}) \vee (h_j(f_i) \neq h_j(f_{i'})) \\ 1 & (f_i \neq f_{i'}) \wedge (h_j(f_i) = h_j(f_{i'})) \end{cases} \quad (2)$$

$I_{i,j,i'} = 1$ iff different flows f_i and $f_{i'}$ are held at the same bucket in the j^{th} array. We use the three situations the same as Section IV-B. We define binary random variable Y_i ($1 \leq i \leq M$) as:

$$Y_i = \begin{cases} 0 & \exists 1 \leq j \leq d, \text{ s.t. } \forall 1 \leq k \leq M, I_{i,j,k} = 0 \\ 1 & \text{else} \end{cases} \quad (3)$$

As mentioned in Subsection III-B, d hash functions $h_1(\dots)h_d(\dots)$ are 2-way independent, and the following proof is based on this condition.

For each flow f_i , if in the d mapped buckets, there is at least one bucket with no hash collision, $Y_i = 0$. Otherwise, in each of these d mapped buckets, \exists a flow f_j ($f_i \neq f_j$) that is also mapped to this bucket, then $Y_i = 1$. So if $Y_i = 0$, for any incoming packet \mathbb{P} belonging to f_i , *Situation 3* can never happen. Now let's calculate $E(Y_i)$, the probability that in each of the d arrays, there are hash collisions in the bucket to which flow f_i is mapped. In a given bucket, the probability that a flow is mapped here is $\frac{1}{w}$, so in a bucket to which f_i is mapped, the probability that no other flow is mapped here is $(1 - \frac{1}{w})^{M-1}$. And in a given array, the probability that hash collision happens in the bucket to which f_i is mapped is $(1 - (1 - \frac{1}{w})^{M-1})$, thus,

$$E(Y_i) = \left[1 - (1 - \frac{1}{w})^{M-1}\right]^d \quad (4)$$

We define random variable $X_{i,j}$ as:

$$X_{i,j} = \sum_{i'=1}^M I_{i,j,i'} n_{i'} Y_i \quad (5)$$

Among the flows held in the same bucket as flow f_i , except for flow f_i itself, some flows are unlikely to cause *Situation 3*, thus unlikely to decay the counter field of this bucket, and others are likely to. $X_{i,j}$ represents the sum of the sizes of the latter kind of flows.

For each incoming packet, if it belongs to flow f_i , the counter field is incremented by 1; if not, the counter field is not changed or decayed. Thus we have

$$n_i - X_{i,j} \leq A_j[h_j(f_i)].C \leq n_i \quad (6)$$

Note that $A_j[h_j(f_i)].C$ is the counter value at the query time. Specifically, if for all packets that do not belong to flow f_i , *Situation 3* happens, and when they are being processed, this counter field is the smallest one in all d mapped buckets, and they all decay the counter field, then $A_j[h_j(f_i)].C = n_i -$

$X_{i,j}$. If all such packets do not decay the counter field, then $A_j[h_j(f_i)].C = n_i$. Then we define random variable $P_{i,j,l}$ as the probability that the l^{th} packet decays the counter field, therefore,

$$A_j[h_j(f_i)].C = n_i - \sum_{l=1}^{X_{i,j}} P_{i,j,l} \quad (7)$$

For any $\epsilon > 0$, we have the following formula based on the Markov inequality.

$$\begin{aligned} & Pr\{A_j[h_j(f_i)].C \leq n_i - \epsilon N\} \\ &= Pr\{n_i - \sum_{l=1}^{X_{i,j}} P_{i,j,l} \leq n_i - \epsilon N\} \\ &= Pr\{\sum_{l=1}^{X_{i,j}} P_{i,j,l} \geq \epsilon N\} \leq \frac{E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})}{\epsilon N} \end{aligned} \quad (8)$$

Now let's focus on $E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$. Recall that in real network traffic, most flows are small, called mouse flows, while a few flows are very large, called elephant flows. Assume that all packets are uniformly distributed. Since we assume that the fingerprint of an elephant flow is held at its mapped bucket since inserted [47], [48], if the l^{th} packet belongs to an elephant flow, *Situation 3* cannot happen at this moment. That is, if the l^{th} packet is to decay the given counter field, it must be a mouse flow and this counter field is the smallest in all d mapped buckets' counter fields.

Recall that $A_j[h_j(f_i)].C$ is the counter value at the query time. We assume that before the query time, when a flow arrives, the counter value is uniformly distributed within the range $[1, A_j[h_j(f_i)].C]$, so the probability that the counter size is equal to any integer within this range is $1/A_j[h_j(f_i)].C$. In addition, the decay happens on condition that 1) the new flow is a mouse flow, whose probability is γ ; 2) *Situation 3* happens and 3) this counter is the first smallest counter. The probability of 2) and 3) is no larger than 1. For any C which satisfies $1 \leq C \leq n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$, we have the following formula:

$$\begin{aligned} Pr\{P_{i,j,l} = \frac{1}{b^C}\} &\leq \frac{\gamma}{A_j[h_j(f_i)].C} \\ &= \frac{\gamma}{n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})} \end{aligned} \quad (9)$$

Let β be $n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$. As a result,

$$\begin{aligned} E(\sum_{l=1}^{X_{i,j}} P_{i,j,l}) &= \sum_{l=1}^{E(X_{i,j})} E(P_{i,j,l}) \\ &\leq E(X_{i,j}) \sum_{C=1}^{\beta} \frac{\gamma}{b^C} \frac{1}{\beta} = \frac{\gamma E(X_{i,j})}{\beta} \cdot \sum_{C=1}^{\beta} \frac{1}{b^C} \\ &= \frac{\gamma E(X_{i,j})}{\beta} \cdot \frac{1}{b} \left[1 - \left(\frac{1}{b}\right)^\beta\right] \\ &\leq \frac{\gamma E(X_{i,j}) \left[1 - \left(\frac{1}{b}\right)^{\beta}\right]}{n_i (b-1)} \end{aligned} \quad (10)$$

Furthermore, for $E(X_{i,j})$, based on Equation 4 and 5,

$$\begin{aligned}
E(X_{i,j}) &= E\left(\sum_{i'=1}^M I_{i,j,i'} n_{i'} Y_i\right) \\
&\leq \sum_{i'=1}^M n_{i'} E(I_{i,j,i'}) E(Y_i) \\
&= \frac{N}{w} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d
\end{aligned} \tag{11}$$

Therefore, based on Equation 10,

$$\begin{aligned}
E\left(\sum_{l=1}^{X_{i,j}} P_{i,j,l}\right) &\leq \frac{\gamma N \left[1 - \left(\frac{1}{b}\right)^{n_i}\right]}{w n_i (b-1)} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d \\
&\leq \frac{\gamma N}{w n_i (b-1)} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d
\end{aligned} \tag{12}$$

Then, based on Equation 8,

$$\begin{aligned}
Pr\{A_j[h_j(f_i)].C \leq n_i - \epsilon N\} &\leq \frac{E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})}{\epsilon N} \\
&\leq \frac{\gamma N}{\epsilon N w n_i (b-1)} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d \\
&= \frac{\gamma}{\epsilon w n_i (b-1)} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d
\end{aligned} \tag{13}$$

For an elephant flow f_i , n_i is very large, so we have

$$\begin{aligned}
Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} &\leq Pr\{\hat{n}_i \leq n_i - \epsilon N\} \\
&\leq \frac{\gamma}{\epsilon w n_i (b-1)} \left[1 - \left(1 - \frac{1}{w}\right)^{M-1}\right]^d \\
&\leq \frac{\gamma}{\epsilon w n_i (b-1)} \left(1 - e^{-\frac{1-M}{w}}\right)^d
\end{aligned}$$

Since w and d are much smaller than M , we have $1 - \delta < (1 - e^{-\frac{1-M}{w}})^d < 1$, where δ is a very small positive number. Therefore, we have

$$Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} \leq \frac{\gamma}{\epsilon w n_i (b-1)}$$

Theorem holds. \square

Theorem 3 is based on an assumption that for an elephant flow, since it is inserted into a bucket, it would be held there all the time. However, if an elephant flow with extremely large size, say 10^{20} , arrives so late that all of its mapped buckets have been filled with other elephant flows with size 1000, it seems impossible to record this flow accurately. This case happens mainly because the current memory size is too small to record elephant flows. Specifically, for an elephant flow f_i , there are the following three situations. 1) This elephant flow f_i arrives early and there are still some empty buckets among its mapped buckets. In this case, f_i is inserted into the empty buckets. f_i can hardly be replaced by other flows due to its high frequency, so in Theorem 3 we assume that such kind of flows are held in the buckets since they are inserted, and we derive mathematical proofs for them in Theorem 3. 2) The elephant flow f_i arrives late but among its d mapped buckets, the smallest counter field is quite small. This means that the flow held in the bucket with the smallest counter field is a

mouse flow, which is easy to be replaced by f_i very soon. After f_i is inserted into this bucket, f_i can hardly be replaced due to its high frequency. Similar to the first case, Theorem 3 can also be applied to this case. 3) The elephant flow f_i arrives late, and all of its d mapped buckets have large counter fields, which means f_i can hardly be inserted into any one of the buckets. Actually, this case typically means the current memory size is too small. Therefore, Theorem 3 only focuses on the first and second cases. For the third case, more memory is needed and we cannot derive any mathematical proofs.

In order to deal with this limitation that elephant flows arriving late are at a disadvantage, we can use the method mentioned in Section III-F. We can use an extra global counter to record how many times a flow's d mapped counters are all large counters. If this extra counter value exceeds the predefined threshold, we add a new array into the HeavyKeeper to make room for the new flow.

In addition, we can observe that in the process of derivation, only $P_{i,j,l}$ is related to the probability decay function. When we choose another decay function, we can derive the formula of $P_{i,j,l}$ in a similar way.

VI. EXPERIMENTAL RESULTS

A. Experiment Setup

Platform: Our experiments are run on a server with dual 6-core CPUs (24 threads, Intel Xeon CPU E5-2620 @2 GHz) and 62 GB total system memory. Each core has two L1 caches with 32KB memory (one instruction cache and one data cache) and one 256KB L2 cache. All cores share one 15MB L3 cache.

Dataset:

1) Campus dataset: This dataset is comprised of IP packets captured from the network of our campus. We rely on the usual definition of a flow, through its 5-tuple, *i.e.*, source IP address, destination IP address, source port, destination port, and protocol type. There are 10M packets in total, belonging to 1M flows.

2) CAIDA dataset: The second dataset is a CAIDA Anonymized Internet Trace from 2016 [49], made of anonymized IP packets. Each flow in this dataset is identified by the source and destination IP address. We use the first 10M. ² packets, belonging to about 4.2M flows.

3) Synthetic datasets: We generate 10 different synthetic datasets according to a Zipf [50] distribution with different skewness (from 0.6 to 3.0)³ Each dataset is comprised of 32M packets, belonging to 1 ~ 10M flows depending on the skewness. Each packet is 4 bytes long. The code of the dataset generator is the one from Web Polygraph [51].

Implementation: The implementation of two versions of HeavyKeeper is done in C++. We also implemented in C++ the other related algorithms including Space-Saving (SS), Lossy

²In network-wide measurement, sketches in different switches are often periodically sent to a collector for timely network traffic analysis. Each period is often small, for example, 10M packets, so we use 10-32M long packet traces.

³Assume there is a stream which has M distinct flows and let N be the total number of flows. Let f_i be the frequency of the i^{th} flow. The skewness γ of this stream refers to $f_i = \frac{N}{i^\gamma \delta(\gamma)}$, where $\delta(\gamma) = \sum_{j=1}^M \frac{1}{j^\gamma}$.

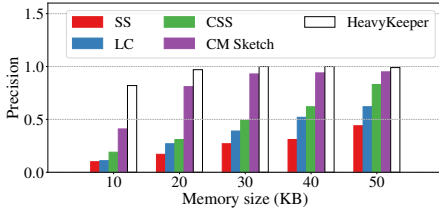


Fig. 4. Precision vs. memory size (Campus).

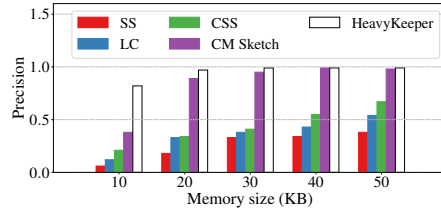


Fig. 5. Precision vs. memory size (CAIDA).

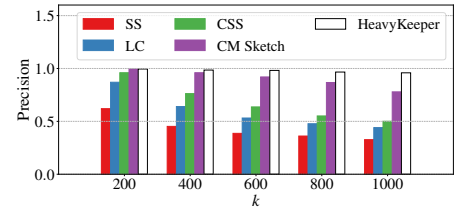


Fig. 6. Precision vs. k (Campus).

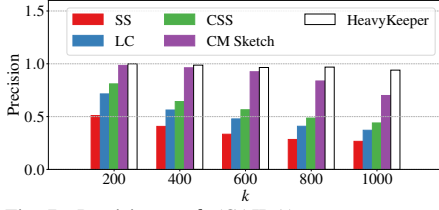


Fig. 7. Precision vs. k (CAIDA).

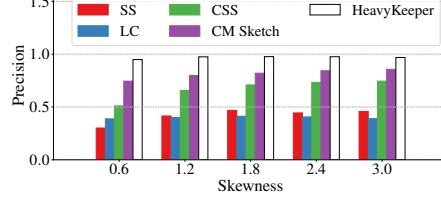


Fig. 8. Precision vs. skewness (Synthetic).

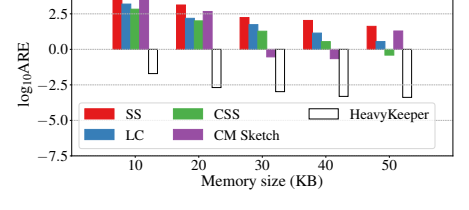


Fig. 9. ARE vs. memory size (Campus).

counting (LC), and the CM sketch⁴ The source code of CSS was provided by its author [31], and is written in Java. It is much slower than Space-Saving written in C++. Therefore, we do not include CSS in our speed experiments. For Space-Saving, Lossy counting, and CSS, the **number of buckets** m is determined by the memory size, which is usually much larger than k . When querying top- k flows, they report the largest k flows of them. For CM sketch, the size of the heap is k , the number of arrays is 3, and the width of each array is determined by the memory size. In our algorithm, the number of buckets m in Stream-Summary is equal to k , and HeavyKeeper occupies the rest memory size. Here we set $d = 2$, and w depends on the memory size. Both the fingerprint field and the counter field are 16-bit long. For experiments on throughput, we ignore operations on the min-heap for the CM sketch, because we can only record flows whose estimated size is larger than a pre-defined threshold.

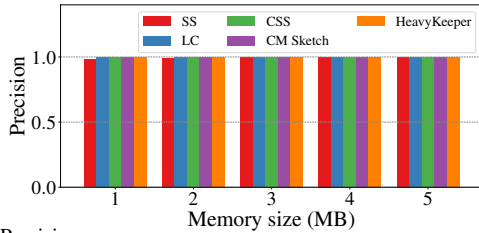


Fig. 10. Precision vs. memory.

B. Metrics

Precision: Precision is defined as $\frac{C}{k}$. Only C flows belong to the real top- k flows.

Average Relative Error (ARE): ARE is defined as $\frac{1}{|\Psi|} \sum_{f_i \in \Psi} \frac{|\hat{n}_i - n_i|}{n_i}$, where Ψ is estimated set of top- k flows, \hat{n}_i is the estimated size of flow f_i , and n_i is the real size of flow f_i . ARE evaluates the error rate reported by the algorithm.

Average Absolute Error (AAE): AAE is defined as $\frac{1}{|\Psi|} \sum_{f_i \in \Psi} |\hat{n}_i - n_i|$, similarly to ARE.

Throughput: We perform insertions of all packets, record the total time used, and calculate the throughput. The throughput

is defined as $\frac{N}{T}$, where N is the total number of packets, and T is the total measured time. We use Million of insertions per second (Mps) to measure the throughput.

C. Experiments on Precision

To achieve a head-to-head comparison, we use the same memory size for each algorithm, and use Hardware Parallel Version as our default version of HeavyKeeper. We perform the experiments for varying memory size and k on the campus and CAIDA datasets, and varying skewness on the synthetic datasets. For experiments of varying memory size, we set $k = 100$, and vary the memory from 10 to 50KB. For experiments of varying k , we set the memory size to 100KB, and vary k from 200 to 1000. For experiments of varying skewness, we set the memory size to 100KB, set $k = 1000$, and vary skewness from 0.6 to 3.0.

Precision vs. memory size: For the campus dataset, when memory size is 10KB (see Figure 4), the precision of Space-Saving, Lossy counting, CSS, and CM sketch is respectively 10%, 11%, 19%, and 41%, while the one of HeavyKeeper is 82%. For the CAIDA dataset (see Figure 5), we find that the precision of HeavyKeeper reaches 99.99% when memory size is larger than 20KB, while for Space-Saving, Lossy counting, CSS, and CM sketch, precision is respectively 18%, 33%, 34%, and 89% when memory size is 50KB.

Precision vs. k : As shown in Figure 6, for the campus dataset, as k becomes larger, the precision of HeavyKeeper stays high, while it degrades for other algorithms. Specifically, the precision of HeavyKeeper is always higher than 95.9%, while that of Space-Saving, Lossy counting, CSS, and CM sketch reaches 32.7%, 44.1%, 50.1%, and 77.9% respectively when $k = 1000$. For the CAIDA dataset (Figure 7), we find that the precision of HeavyKeeper is always above 94%, while for Space-Saving, Lossy counting, CSS, and CM sketch, it is 26.6%, 37.1%, 44%, and 70% respectively when $k = 1000$.

Precision vs. skewness: As shown in Figure 8, the precision of HeavyKeeper reaches 99.99%. For all considered values of skewness, the precision of HeavyKeeper does not go below 94.9%, while the highest precision for Space-Saving, Lossy counting, CSS, and CM sketch is 46.8%, 41.3%, 74.5%, and 85.7%, respectively.

⁴There is an open source library [52] that implements Lossy Counting, the CM Sketch, Space Saving, and others. Because the format of packets is different from our datasets, we implemented these algorithms by ourselves.

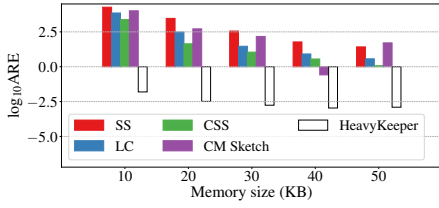


Fig. 11. ARE vs. memory size (CAIDA).

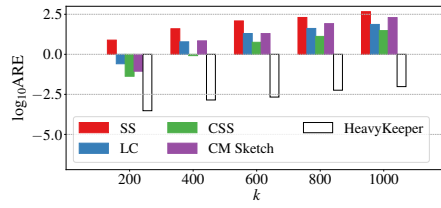


Fig. 12. ARE vs. k (Campus).

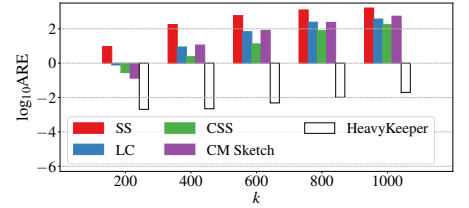


Fig. 13. ARE vs. k (CAIDA).

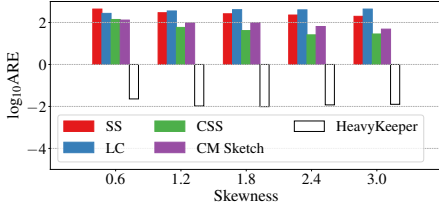


Fig. 14. ARE vs. skewness (Synthetic).

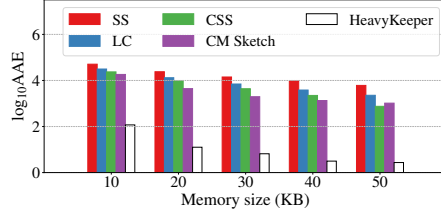


Fig. 15. AAE vs. memory size (Campus).

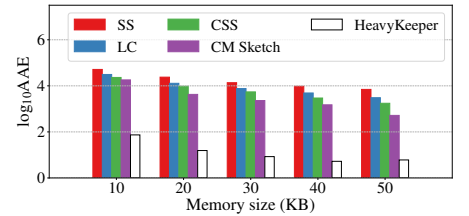


Fig. 16. AAE vs. memory size (CAIDA).

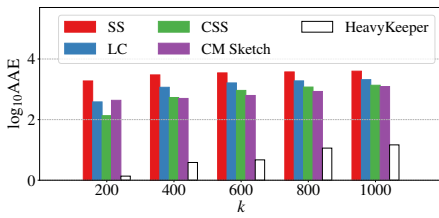


Fig. 17. AAE vs. k (Campus).

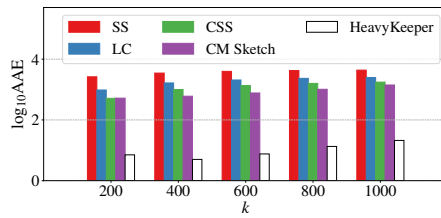


Fig. 18. AAE vs. k (CAIDA).

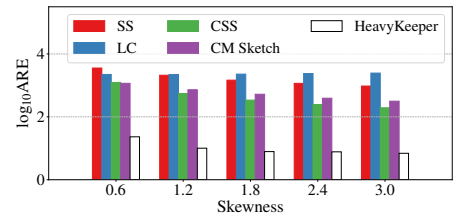


Fig. 19. AAE vs. skewness (Synthetic).

D. Experiments on AAE and ARE

In this section, we focus on the ARE and the AAE of the estimated frequency of reported top- k flows. We also conduct experiments with varying memory size, k , and skewness. The parameter settings are the same as in Section VI-C.

ARE vs. memory size: As shown in Figure 9, for the campus dataset, we find that the ARE of HeavyKeeper is smaller than 0.01 when memory size is larger than 20KB, while for Space-Saving, Lossy counting, CSS, and CM sketch, it is larger than 100. For the CAIDA dataset (see Figure 11), we find that the ARE of HeavyKeeper is between 21119 and 1190365 times smaller than the one of Space-Saving, between 2955 and 456275 times smaller than the one of Lossy counting, between 950 and 154047 times smaller than the one of CSS, and between 238 and 656145 times smaller than the one of CM sketch.

ARE vs. k : As shown in Figure 12, for the campus dataset, we find that the ARE of HeavyKeeper is between 25579 and 56791 times smaller than the one of Space-Saving, between 852 and 9312 times smaller than the one of Lossy counting, between 142 and 3132 times smaller than the one of CSS, and between 293 and 20370 times smaller than the one of CM sketch. For the CAIDA dataset (see Figure 13), we find that the ARE of HeavyKeeper is between 4506 and 121912 times smaller than the one of Space-Saving, between 383 and 23666 times smaller than the one of Lossy counting, between 137 and 8816 times smaller than the one of CSS, and between 66 and 27290 times smaller than the one of CM sketch.

ARE vs. skewness: As shown in Figure 14, for all considered values of skewness, we find that the ARE of HeavyKeeper is between 15566 and 27829 times smaller than that of Space-Saving, between 11915 and 41575 times smaller than that of Lossy counting, between 2174 and 6099 times smaller than

that of CSS, and between 3819 and 10080 times smaller than that of CM sketch.

AAE vs. memory size: As shown in Figure 15, for the campus dataset, we find that the AAE of HeavyKeeper is between 433 and 3013 times smaller than that of Space-Saving, between 267 and 1221 times smaller than that of Lossy counting, between 200 and 758 times smaller than that of CSS, and between 155 and 428 times smaller than that of CM sketch. For the CAIDA dataset (see Figure 16), we find that the AAE of HeavyKeeper is between 697 and 1810 times smaller than that of Space-Saving, between 421 and 928 times smaller than that of Lossy counting, between 289 and 647 times smaller than the one of CSS, and between 86 and 284 times smaller than that of CM sketch.

AAE vs. k : As shown in Figure 17, for the campus dataset, we find that the AAE of HeavyKeeper is between 271 and 1382 times smaller than that of Space-Saving, between 142 and 346 times smaller than that of Lossy counting, between 93 and 196 times smaller than that of CSS, and between 74 and 318 times smaller than that of CM sketch. For CAIDA dataset (see Figure 18), we find that the AAE of HeavyKeeper is between 206 and 694 times smaller than that of Space-Saving, between 118 and 329 times smaller than that of Lossy counting, between 73 and 199 times smaller than that of CSS, and between 67 and 121 times smaller than that of CM sketch.

AAE vs. skewness: From Figure 19, we find that the AAE of HeavyKeeper is between 137 and 209 times smaller than that of Space-Saving, between 96 and 355 times smaller than that of Lossy counting, between 28 and 55 times smaller than that of CSS, and between 45 and 73 times smaller than that of CM sketch.

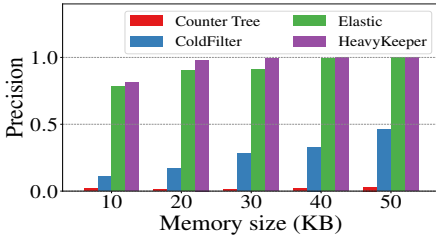


Fig. 20. Precision vs. memory size.

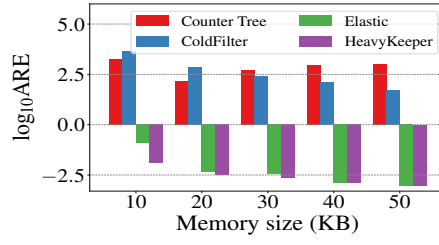


Fig. 21. ARE vs. memory size.

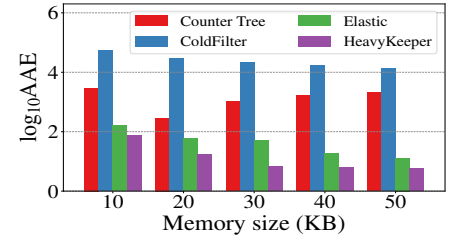


Fig. 22. AAE vs. memory size.

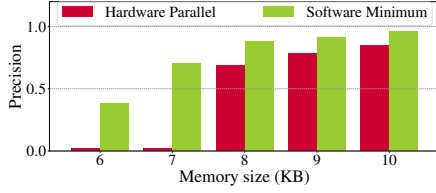


Fig. 23. Precision vs. memory size.

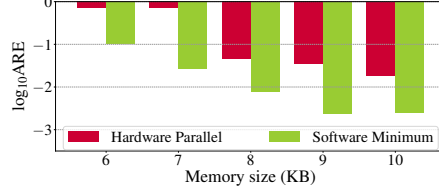


Fig. 24. ARE vs. memory size.

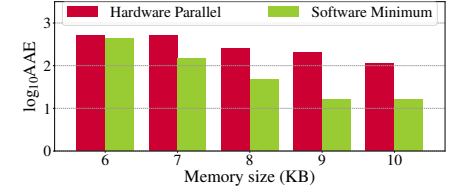


Fig. 25. AAE vs. memory size.

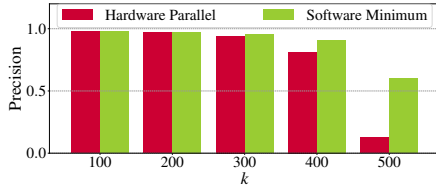


Fig. 26. Precision vs. k .

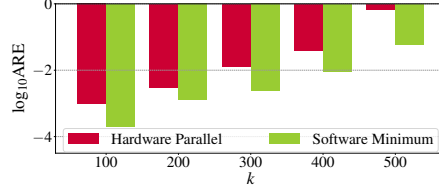


Fig. 27. ARE vs. k .

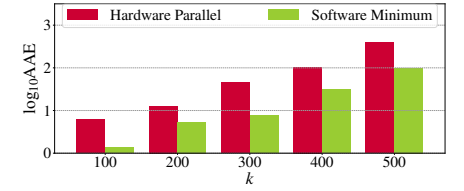


Fig. 28. AAE vs. k .

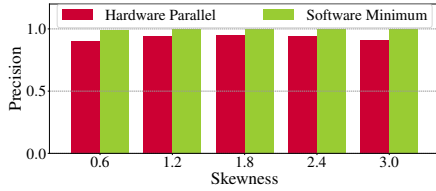


Fig. 29. Precision vs. skewness.

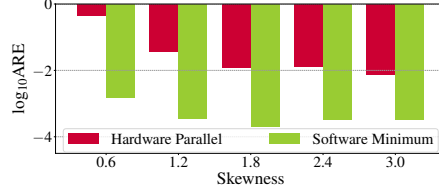


Fig. 30. ARE vs. skewness.

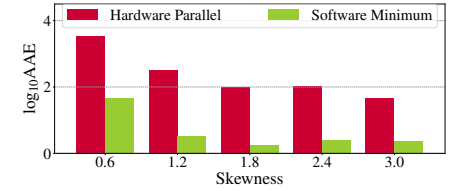


Fig. 31. AAE vs. skewness.

E. Compare with Recent Works

In this section, we compare our algorithm with recent works. First we show the differences between HeavyKeeper and HeavyGuardian. Then we compare our HeavyKeeper with the Elastic sketch, Counter Tree and Cold Filter. For the Elastic sketch and Cold Filter, the source codes are from their authors [41], [43]. We use Cold Filter with Space Saving to evaluate its performance, because the performance of Cold Filter with Space Saving is the best in that paper. For Counter Tree, we use the formulas derived from its author [44] to estimate frequencies of flows. We only report results for the campus dataset by varying the memory size. Here we set $k = 100$ and vary memory size from 10KB to 50KB.

As mentioned before in Section I-A, HeavyGuardian can also find items with large frequencies, but we do not compare our HeavyKeeper with HeavyGuardian, due to the following three differences. 1) These two algorithms have different focuses. HeavyGuardian focuses on generality. It can handle five different tasks: frequency estimation, heavy hitter detection, heavy change detection, frequency distribution estimation, and entropy estimation. But it was not applied to find top- k elephant flows. Our HeavyKeeper is designed to only find top- k elephant flows accurately. 2) HeavyGuardian is the first algorithm that supports real-time entropy estimation, but HeavyKeeper cannot handle real-time entropy estimation. 3)

HeavyGuardian has the above advantages at the cost of being applicable for software platforms only, *i.e.*, it cannot be implemented on hardware platforms. While in our HeavyKeeper for Hardware Parallel version, the operation in each array can be implemented in parallel on hardware platforms. Therefore, we do not compare our algorithm with HeavyGuardian. We compare HeavyKeeper with the Elastic sketch, Counter Tree and Cold Filter, which is detailed as follows.

Measuring precision: As shown in Figure 20, the precision of HeavyKeeper is much better than Counter Tree and Cold Filter. Next we explain the reason of the performance difference between our algorithm and others. For Counter Tree, it uses formulas to estimate frequencies of flows, which might cause large error. For Cold Filter, its key data structure is Space Saving [29], whose performance is worse than HeavyKeeper, and the cold filter takes up a certain amount of memory. For the Elastic sketch, it is a general data structures, while HeavyKeeper just focuses on finding top- k elephant flows. That is why HeavyKeeper is slightly better than the Elastic sketch.

Measuring ARE: As shown in Figure 21, the ARE of HeavyKeeper is the smallest compared with other recent works. Specifically, when the memory size is 10KB, the ARE of Counter Tree, Cold Filter and the Elastic sketch are $10^{3.2}$, $10^{3.6}$ and $10^{-0.9}$, respectively, while that of HeavyKeeper

is smaller than $10^{-1.8}$. This indicates HeavyKeeper could handle the situation in tight memory much better than other algorithms.

Measuring AAE: As shown in Figure 22, the AAE of HeavyKeeper is the smallest compared with other recent works. Specifically, when the memory size is 10KB, the AAE of Counter Tree, Cold Filter and the Elastic sketch are $10^{3.4}$, 10^4 and $10^{2.1}$, respectively, while that of HeavyKeeper is smaller than $10^{1.9}$. As the memory size increases, the AAE of our algorithm is always the smallest compared with other algorithms.

F. Performance on Very Big Dataset

We also conduct experiments on very big datasets. We set $k = 1000$ and the memory size to 100KB. For every 10M packets, we report top- k elephant flows and evaluate the precision by comparing with real top- k elephant flows. As shown in Figure 32, as the total number of packets increases, the precision slightly reduces. However, we can observe that the precision still reaches a high value when the total number of packets is 10^8 .

G. Hardware Parallel Version vs. Software Minimum Version

In this section, we compare Hardware Parallel Version with Software Minimum Version. We conduct experiments with varying memory size, k , and skewness. Due to the high accuracy of our algorithm, we set the smaller memory size to show the difference of performance between two versions clearly. Specifically, for experiments of varying memory size, we set $k = 100$, and vary the memory size from 6KB to 10KB; for experiments of varying k , we set the memory size to 30KB, and vary k from 100 to 500; for experiments of varying skewness, we set the memory size to 10KB and $k = 100$. Since the results are similar on CAIDA and campus datasets, we just show the performance of two versions on campus dataset.

Varying memory size: As shown in Figure 23, when memory size is 5KB or 6KB, the precision of Hardware Parallel Version is only 2%, and the reason behind is that there are only a few buckets, which cannot record all the largest k flows. On the other hand, the precision of Software Minimum Version achieves 38% and 70% when memory size is 5KB and 6KB, respectively, and the reason behind is that each flow has no duplicate when it is inserted into the hash table, and therefore the Software Minimum Version saves memory more efficiently. As shown in Figure 24 and 25, we find that the ARE and AAE of Software Minimum Version are smaller than those of Hardware Parallel Version.

Varying k : As shown in Figure 26, as k increases, the precision of Hardware Parallel Version decreases from 100% to 13%, while the Software Minimum Version still achieves 60% precision when $k = 1000$. As shown in Figure 27 and 28, we find that the ARE and AAE of Software Minimum Version are smaller than those of Hardware Parallel Version.

Varying skewness: As shown in Figure 29-31, for all considered values of skewness, the precision of Software Minimum Version is always larger than that of Hardware Parallel Version, and the ARE and AAE of Software Minimum Version are always smaller than those of Hardware Parallel Version.

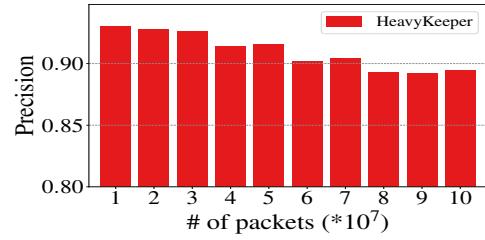


Fig. 32. Precision vs. # of packets.

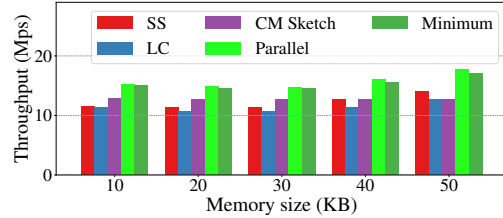


Fig. 33. Throughput vs. memory size.

H. Experiments on Throughput

We now turn to the throughput of the algorithms. We only report results for the campus dataset due to space limitations. We set $k = 100$, and vary memory size from 10KB to 50KB. Since our server of running experiments is much older than most of the current ones, the throughput of experimental results might be slightly lower than the results in other papers.

Throughput vs. memory size: As shown in Figure 33, we find that the throughput of HeavyKeeper is always higher than other algorithms, and the throughput of HeavyKeeper of Hardware Parallel Version is slightly higher than the Software Minimum Version. Indeed, the average throughput of HeavyKeeper of Hardware Parallel Version and Software Minimum Version is 15.52Mps, 15.27Mps, respectively, while it is 12.15Mps, 11.34Mps, and 12.72Mps for Space-Saving, Lossy counting, and CM sketch. These results show that HeavyKeeper not only is more accurate than previous work, but also achieves higher throughput as well.

VII. OPEN vSWITCH DEPLOYMENT

In this section, we implement our HeavyKeeper algorithm on a software switch platform: Open vSwitch (OVS). We will present details of our implementation and experimental results to show the performance running on Open vSwitch.

A. OVS Implementation

The OVS implementation of our HeavyKeeper algorithm consists of three components: 1) the modified OVS datapath, 2) the shared memory buffering flow IDs, and 3) the user-space program of HeavyKeeper processing flow IDs. For each incoming packet, it will be first inserted into the OVS datapath for forwarding. Besides, we modify the source codes of OVS datapath to parse the flow ID of the incoming packet, and then insert its flow ID into the shared memory (the shared memory is created initially). Finally, the user-space program will read the flow IDs from the shared memory, and process them as incoming packets.

B. OVS Evaluation

We use synthetic trace to conduct experiments in OVS with 4 threads and 40G link min-size packets to evaluate the throughput of HeavyKeeper and other algorithms. In order

to improve the performance of OVS, we integrate OVS with DPDK (Data Plane Development Kit). DPDK implements the datapath entirely in the user-space, and thus it eliminates the overhead of a context switch and memory copies between user-space and kernel-space. Besides, we also show the throughput of OVS without using any algorithm to show the impact of algorithms. We set the memory size to 50KB.

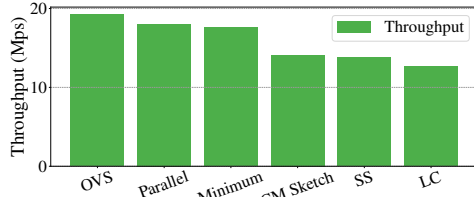


Fig. 34. Throughput on OVS platform.

As shown in Figure 34, the throughput of HeavyKeeper is near the original throughput of OVS. Specifically, the throughput of the original OVS is 19.22Mbps, and that of HeavyKeeper of Hardware Parallel Version and Software Minimum Version is 18.03Mbps, 17.62Mbps, respectively. However, the throughput of CM sketch, Space-Saving, and Lossy Counting is 14.14Mbps, 13.80Mbps, and 12.64Mbps, respectively. The results show that our HeavyKeeper algorithm has little impact to the performance of OVS, while other algorithms decrease the throughput significantly.

VIII. CONCLUSION

Finding the top- k elephant flows is a critical task for network traffic measurement. Existing algorithms for finding top- k flows cannot achieve high precision when traffic speed is high and memory usage is small. In this paper, we propose a novel data structure, called HeavyKeeper, which achieves a much higher precision on top- k queries and a much lower error rate on flow size estimation, compared to previous algorithms. The key idea of HeavyKeeper is that it intelligently omits mouse flows, and focuses on recording the information of elephant flows by using the exponential-weakening decay strategy. Our evaluation confirms that HeavyKeeper achieves 99.99% precision for finding the top- k elephant flows, while also achieving a reduction in the error rate of the estimated flow size by about 3 orders of magnitude compared to the state-of-the-art algorithms. We have released the source code of HeavyKeeper and all related algorithms at GitHub [46].

ACKNOWLEDGMENT

We would like to thank the anonymous reviewers for their insightful suggestions. This work is supported by Primary Research & Development Plan of China (2018YFB1004403, 2016YFB1000304), NSFC (61672061), and Shenzhen Peacock Innovation Program (KQJSCX20180323174744219).

REFERENCES

[1] J. Gong, T. Yang, H. Zhang, H. Li, S. Uhlig, S. Chen, L. Uden, and X. Li, "Heavykeeper: An accurate algorithm for finding top-k elephant flows," in *Proc. USENIX ATC*, 2018, pp. 909–921.
 [2] A. Sivaraman, S. Subramanian, and et al., "Programmable packet scheduling at line rate," in *Proc. ACM SIGCOMM*, 2016.

[3] A. Feldmann, A. Greenberg, and et al., "Deriving traffic demands for operational ip networks: Methodology and experience," in *Proc. ACM SIGCOMM*, 2000.
 [4] A. Lakhina, M. Crovella, and C. Diot, "Characterization of network-wide anomalies in traffic flows," in *Proc. ACM IMC*, 2004.
 [5] O. Rottenstreich and J. Topolcai, "Optimal rule caching and lossy compression for longest prefix matching," *IEEE/ACM Transactions on Networking*, vol. 25, no. 2, pp. 864–878, 2017.
 [6] Q. Huang, P. P. Lee, and Y. Bao, "Sketchlearn: relieving user burdens in approximate measurement with automated statistical inference," in *ACM SIGCOMM 2018*, pp. 576–590.
 [7] Q. Huang, X. Jin, P. P. Lee, R. Li, L. Tang, Y.-C. Chen, and G. Zhang, "Sketchvisor: Robust network measurement for software packet processing," in *ACM SIGCOMM 2017*, pp. 113–126.
 [8] L. Tang, Q. Huang, and P. P. C. Lee, "Mv-sketch: A fast and compact invertible sketch for heavy flow detection in network data streams," in *IEEE INFOCOM 2019*.
 [9] C. Hu, B. Liu, H. Zhao, and et al., "Disco: Memory efficient and accurate flow statistics for network measurement," in *IEEE ICDCS*, 2010, pp. 665–674.
 [10] C. Hu, B. Liu, H. Zhao, K. Chen, Y. Chen, C. Wu, and Y. Cheng, "Disco: Memory efficient and accurate flow statistics for network measurement," in *Distributed Computing Systems (ICDCS), 2010 IEEE 30th International Conference on*. IEEE, 2010, pp. 665–674.
 [11] H. C. Zhao, A. Lall, M. Ogihara, O. Spatscheck, J. Wang, and J. Xu, "A data streaming algorithm for estimating entropies of od flows," in *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*. ACM, 2007, pp. 279–290.
 [12] H. Dai, M. Shahzad, A. X. Liu, and Y. Zhong, "Finding persistent items in data streams," *Proceedings of the VLDB Endowment*, vol. 10, no. 4, pp. 289–300, 2016.
 [13] A. Kumar, J. Xu, and J. Wang, "Space-code bloom filter for efficient per-flow traffic measurement," *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 12, pp. 2327–2339, 2006.
 [14] O. Rottenstreich, Y. Kanizo, and I. Keslassy, "The variable-increment counting bloom filter," *IEEE/ACM Transactions on Networking (TON)*, vol. 22, no. 4, pp. 1092–1105, 2014.
 [15] S. Z. Kiss, E. Hosszu, J. Topolcai, L. Rónyai, and O. Rottenstreich, "Bloom filter with a false positive free zone," *Proc. IEEE INFOCOM, Honolulu, HI, USA*, 2018.
 [16] K. Mirylenka, G. Cormode, T. Palpanas, and D. Srivastava, "Conditional heavy hitters: detecting interesting correlations in data streams," *The VLDB Journal*, vol. 24, no. 3, pp. 395–414, 2015.
 [17] J. H. Chang and W. S. Lee, "Finding recent frequent itemsets adaptively over online data streams," in *Proc. ACM SIGKDD*. ACM, 2003, pp. 487–492.
 [18] Y.-L. Cheung and A. W.-C. Fu, "Mining frequent itemsets without support threshold: with and without item constraints," *IEEE TKDE*, vol. 16, no. 9, pp. 1052–1069, 2004.
 [19] G. Salton and M. J. McGill, "Introduction to modern information retrieval," 1986.
 [20] M. A. Soliman, I. F. Ilyas, and K. C.-C. Chang, "Top-k query processing in uncertain databases," in *Proc. IEEE ICDE*, 2007, pp. 896–905.
 [21] Y. Zhang, B. Fang, and Y. Zhang, "Identifying heavy hitters in high-speed network monitoring," *Science China Information Sciences*, vol. 53, no. 3, pp. 659–676, 2010.
 [22] P. Roy, A. Khan, and G. Alonso, "Augmented sketch: Faster and more accurate stream processing," in *Proc. SIGMOD 2016*.
 [23] G. Cormode, "Sketch techniques for approximate query processing," *Foundations and Trends in Databases. NOW publishers*, 2011.
 [24] G. Cormode and S. Muthukrishnan, "An improved data stream summary: the count-min sketch and its applications," *Journal of Algorithms*, vol. 55, no. 1, pp. 58–75, 2005.
 [25] C. Estan and G. Varghese, *New directions in traffic measurement and accounting*. ACM, 2002, vol. 32, no. 4.
 [26] Y. Zhang, M. Roughan, W. Willinger, and L. Qiu, "Spatio-temporal compressive sensing and internet traffic matrices," in *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 4. ACM, 2009, pp. 267–278.
 [27] T. Benson, A. Akella, and D. A. Maltz, "Network traffic characteristics of data centers in the wild," in *Proceedings of the 10th ACM SIGCOMM conference on Internet measurement*. ACM, 2010, pp. 267–280.
 [28] D. Maltz, "Unraveling the complexity of network management," 2009.
 [29] Z. Li, F. Xiao, S. Wang, T. Pei, and J. Li, "Achievable rate maximization for cognitive hybrid satellite-terrestrial networks with af-relays," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 2, pp. 304–313, 2018.

- [30] M. H. ur Rehman, C. S. Liew, A. Abbas, P. P. Jayaraman, T. Y. Wah, and S. U. Khan, "Big data reduction methods: a survey," *Data Science and Engineering*, vol. 1, no. 4, pp. 265–284, 2016.
- [31] R. Ben-Basat, G. Einziger, R. Friedman, and Y. Kassner, "Heavy hitters in streams and sliding windows," in *Proc. IEEE INFOCOM*, 2016.
- [32] A. Metwally, D. Agrawal, and A. El Abbadi, "Efficient computation of frequent and top-k elements in data streams," in *Proc. Springer ICDT 2005*.
- [33] M. Charikar, K. Chen, and M. Farach-Colton, "Finding frequent items in data streams," *Automata, languages and programming*, pp. 784–784, 2002.
- [34] G. S. Manku and R. Motwani, "Approximate frequency counts over data streams," in *Proc. VLDB*, 2002, pp. 346–357.
- [35] N. B. Seghouani, F. Bugiotti, M. Hewasinghage, S. Isaj, and G. Quercini, "A frequent named entities-based approach for interpreting reputation in twitter," *Data Science and Engineering*, vol. 3, no. 2, pp. 86–100, 2018.
- [36] K. Li and G. Li, "Approximate query processing: what is new and where to go?" *Data Science and Engineering*, vol. 3, no. 4, pp. 379–397, 2018.
- [37] Z. Li, B. Chang, S. Wang, A. Liu, F. Zeng, and G. Luo, "Dynamic compressive wide-band spectrum sensing based on channel energy reconstruction in cognitive internet of things," *IEEE Transactions on Industrial Informatics*, 2018.
- [38] Y. Li, R. Miao, C. Kim, and M. Yu, "Flowradar: A better netflow for data centers," in *NSDI*, 2016, pp. 311–324.
- [39] W. Feng and H. Mounir, "Matching the speed gap between sram and dram," in *Proc. IEEE HSPR*, 2008, pp. 104–109.
- [40] E. Demaine, A. López-Ortiz, and J. Munro, "Frequency estimation of internet packet streams with limited space," *AlgorithmsESA 2002*, pp. 11–20, 2002.
- [41] T. Yang, J. Jiang, P. Liu, Q. Huang, J. Gong, Y. Zhou, R. Miao, X. Li, and S. Uhlig, "Elastic sketch: Adaptive and fast network-wide measurements," in *Proceedings of the 2018 Conference of the ACM SIGCOMM*. ACM, 2018, pp. 561–575.
- [42] T. Yang, J. Gong, H. Zhang, L. Zou, L. Shi, and X. Li, "Heavyguardian: Separate and guard hot items in data streams," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2018, pp. 2584–2593.
- [43] Y. Zhou, T. Yang, J. Jiang, B. Cui, M. Yu, X. Li, and S. Uhlig, "Cold filter: A meta-framework for faster and more accurate stream processing," in *Proceedings of the 2018 International Conference on Management of Data*. ACM, 2018, pp. 741–756.
- [44] C. Min and S. Chen, "Counter tree: A scalable counter architecture for per-flow traffic measurement," *IEEE/ACM Transactions on Networking*, vol. PP, no. 99, pp. 1–14, 2017.
- [45] G. Einziger and R. Friedman, "Counting with tinytable: Every bit counts!" in *Proc. ICDCN 2016*.
- [46] "The source codes of heavykeeper and other related algorithms," <https://github.com/papergitkeeper/heavy-keeper-project>.
- [47] S. Balaji, M. N. Krishnan, M. Vajha, V. Ramkumar, B. Sasidharan, and P. V. Kumar, "Erasure coding for distributed storage: an overview," *Science China Information Sciences*, vol. 61, no. 10, p. 100301, 2018.
- [48] X. Tang, S.-T. Xia, C. Tian, Q. Huang, and X.-G. Xia, "Special focus on distributed storage coding," *Science China Information Sciences*, vol. 61, no. 10, p. 100300, 2018.
- [49] "The caida anonymized internet traces 2016," <http://www.caida.org/data/overview/>.
- [50] D. M. Powers, "Applications and explanations of zipf's law," in *Proceedings of the joint conferences on new methods in language processing and computational natural language learning*. Association for Computational Linguistics, 1998, pp. 151–160.
- [51] A. Rousskov and D. Wessels, "High-performance benchmarking with web polygraph," *Software: Practice and Experience*, 2004.
- [52] "The open source library that implements lossy counting, cm sketch, space saving, and others," <http://hadjieleftheriou.com/frequent-items>.



Tong Yang received his PHD degree in Computer Science from Tsinghua University in 2013. He visited Institute of Computing Technology, Chinese Academy of Sciences (CAS). Now he is a research assistant professor in Computer Science Department, Peking University. His research interests include network measurements, sketches, IP lookups, Bloom filters, sketches and KV stores. He published papers in SIGCOMM, SIGKDD, SIGMOD, SIGCOMM CCR, VLDB, ATC, ToN, ICDE, INFOCOM, etc.



Haowei Zhang is a sophomore student in Peking University, advised by Tong Yang. His research interests include network measurement and data stream processing systems. He has some publications in the area of networking and data stream processing.



Jinyang Li is a sophomore student in Peking University, advised by Tong Yang. Her research interests include network measurements, sketches, Bloom filters, data stream processing, and hash tables.

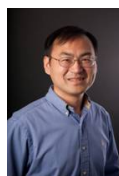


Junzhi Gong is a senior student in Peking University, advised by Tong Yang. His research interests include network measurement and data stream processing systems. He has some publications in the area of networking and data stream processing.



Steve Uhlig obtained a Ph.D. degree in Applied Sciences from the University of Louvain, Belgium, in 2004. From 2004 to 2006, he was a Postdoctoral Fellow of the Belgian National Fund for Scientific Research (F.N.R.S.). His thesis won the annual IBM Belgium/F.N.R.S. Computer Science Prize 2005. Between 2004 and 2006, he was a visiting scientist at Intel Research Cambridge, UK, and at the Applied Mathematics Department of University of Adelaide, Australia. Between 2006 and 2008, he was with Delft University of Technology, the Netherlands.

Prior to joining Queen Mary, he was a Senior Research Scientist with Technische Universitt Berlin/Deutsche Telekom Laboratories, Berlin, Germany. Starting in January 2012, he is the Professor of Networks and Head of the Networks Research group at Queen Mary, University of London. Between 2012 and 2016, he was a guest professor at the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China.



Shigang Chen received the Ph.D. degrees in computer science from the University of Illinois at 1999. He is currently a Professor with the Department of Computer and Information Science and Engineering, University of Florida. He has authored over 160 peer-reviewed journal/conference papers. His research interests include computer networks, Internet security, wireless communications, and distributed computing. He served in various chair positions or as committee members for numerous conferences. He is an ACM Distinguished Member and a Distinguished Lecturer of the IEEE Communication Society. He was an Associate Editor for the IEEE/ACM TRANSACTIONS ON NETWORKING, the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and a number of other journals.



Xiaoming Li is a professor in computer science and technology and the director of Institute of Network Computing and Information Systems (NCIS) at Peking University, China. His current research interest is in search engine and web mining. He led the effort of developing a Chinese search engine (Tianwang) since 1999, and is the founder of the Chinese web archive (Web InfoMall).

APPENDIX

In this appendix, we first present the pseudo-codes of HeavyKeeper in the Parallel version and Minimum version, second we prove that there is no over-estimation, and finally we derive the formula of error bound of HeavyKeeper in the basic version.

A. Pseudo-codes

Algorithm 1 and 2 show the insertion process of the Parallel version and Minimum version, respectively.

Algorithm 1: Insertion process for finding top- k flows (Parallel version).

Input: A packet \mathbb{P}_l belonging to flow f_i

```

1 flag ← false;
2 if  $f_i \in \text{min\_heap}$  then
3   flag ← true;
4 HeavyK_V ← 0;
5 for  $j \leftarrow 1$  to  $d$  do
6    $C \leftarrow A_j[h_j(f_i)].C$ ;
7   if  $C = 0$  then
8      $A_j[h_j(f_i)].FP \leftarrow \mathbb{F}_i$ ;
9      $A_j[h_j(f_i)].C \leftarrow 1$ ;
10    HeavyK_V ← max(HeavyK_V, 1);
11  else if  $A_j[h_j(f_i)].FP = \mathbb{F}_i$  then
12    if flag = true or  $C < \text{min\_heap}.n_{\text{min}}$  then
13       $A_j[h_j(f_i)].C ++$ ;
14      HeavyK_V ←
15        max(HeavyK_V,  $A_j[h_j(f_i)].C$ );
16    else if rand() <  $b^{-C}$  then
17       $A_j[h_j(f_i)].C --$ ;
18      if  $A_j[h_j(f_i)].C = 0$  then
19         $A_j[h_j(f_i)].FP \leftarrow \mathbb{F}_i$ ;
20         $A_j[h_j(f_i)].C \leftarrow 1$ ;
21        HeavyK_V ← max(HeavyK_V, 1);
22  if flag = true then
23    min_heap[ $f_i$ ] ← max(HeavyK_V, min_heap[ $f_i$ ]);
24  else
25    if min_heap has empty buckets or
26      HeavyK_V -  $n_{\text{min}} = 1$  then
27      min_heap.insert( $f_i$ );

```

B. Proof of No Over-estimation Error of HeavyKeeper

Theorem 4. Let $n_i(t)$ be the real size of flow f_i at time t , and let $A_j[h_j(f_i)](t).C$ be the counter field of the mapped bucket of flow f_i in the j^{th} array at time t . If there is no fingerprint collision, then

$$\forall j, t, A_j[h_j(f_i)](t).C \leq n_i(t) \quad (14)$$

Proof. When $t = 0$, no packet maps into this bucket, so $n_i(0) = 0$ and $A_j[h_j(f_i)](0).C = 0$. Therefore, the theorem holds at time 0. Let's now prove by induction that the theorem holds at any time.

Algorithm 2: Insertion process for finding top- k flows (Minimum version).

Input: A packet \mathbb{P}_l belonging to flow f_i

```

1 flag ← false;
2 if  $f_i \in \text{min\_heap}$  then
3   flag ← true;
4 min_count ←  $A_1[h_1(f_i)].C$ ;
5 min_array ← 1;
6 first_empty ← 0;
7 HeavyK_V ← 0;
8 add_flag ← false;
9 for  $j \leftarrow 1$  to  $d$  do
10  if  $A_j[h_j(f_i)].FP = \mathbb{F}_i$  then
11    if flag = true or  $A_j[h_j(f_i)].C <
12      \text{min\_heap}.n_{\text{min}}$  then
13       $A_j[h_j(f_i)].C ++$ ;
14      add_flag ← true;
15      HeavyK_V ←  $A_j[h_j(f_i)].C$ ;
16      break;
17  else
18    if  $A_j[h_j(f_i)].FP = \text{null}$  and first_empty = 0
19    then
20      first_empty ←  $j$ ;
21    else
22      if  $A_j[h_j(f_i)].FP \neq \text{null}$  and
23         $A_j[h_j(f_i)].C < \text{min\_count}$  then
24        min_count ←  $A_j[h_j(f_i)].C$ ;
25        min_array ←  $j$ ;
26  if add_flag = false then
27    if first_empty > 0 then
28       $A_{\text{first\_empty}}[h_{\text{first\_empty}}(f_i)].FP \leftarrow \mathbb{F}_i$ ;
29       $A_{\text{first\_empty}}[h_{\text{first\_empty}}(f_i)].C \leftarrow 1$ ;
30      HeavyK_V ← 1;
31    else
32      if rand() <  $b^{-C}$  then
33         $A_{\text{min\_array}}[h_{\text{min\_array}}(f_i)].C --$ ;
34        if  $A_{\text{min\_array}}[h_{\text{min\_array}}(f_i)].C = 0$  then
35           $A_{\text{min\_array}}[h_{\text{min\_array}}(f_i)].FP \leftarrow \mathbb{F}_i$ ;
36           $A_{\text{min\_array}}[h_{\text{min\_array}}(f_i)].C \leftarrow 1$ ;
37          HeavyK_V ← 1;
38  if flag = true then
39    min_heap[ $f_i$ ] ← max(HeavyK_V, min_heap[ $f_i$ ]);
40  else
41    if min_heap has empty buckets or
42      HeavyK_V -  $n_{\text{min}} = 1$  then
43      min_heap.insert( $f_i$ );

```

When $t = 0$, the theorem holds.

If the theorem holds when $t = v$, let's prove that the theorem also holds when $t = v + 1$.

For HeavyKeeper in the basic version, there are three cases when $t = v + 1$:

Case 1: The new incoming packet is not mapped to bucket $A_j[h_j(f_i)]$. Then $n_i(v+1) = n_i(v)$ and $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C$. Therefore, $A_j[h_j(f_i)](v+1).C \leq n_i(v+1)$.

Case 2: The new incoming packet belongs to flow f_i . Then $n_i(v+1) = n_i(v) + 1$ and $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C + 1$. Therefore, $A_j[h_j(f_i)](v+1).C \leq n_i(v+1)$.

Case 3: The new incoming packet is mapped to bucket $A_j[h_j(f_i)]$ but does not belong to flow f_i . Then $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C$ or $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C - 1$, and $n_i(v+1) = n_i(v)$. Therefore, $A_j[h_j(f_i)](v+1).C \leq n_i(v+1)$.

Therefore, for HeavyKeeper in the basic version, for any time t ,

$$A_j[h_j(f_i)](t).C \leq n_i(t)$$

For HeavyKeeper in the minimum version, Case 1 and Case 2 are the same. We just focus on Case 3.

Case 3 in the Minimum Version: The new incoming packet is mapped to bucket $A_j[h_j(f_i)]$ but does not belong to flow f_i . Suppose the new packet belongs to flow f_j ($f_i \neq f_j$). Among the other $d-1$ mapped buckets, if there is an empty bucket or there is a bucket holding the fingerprint \mathbb{F}_j , then HeavyKeeper does nothing to this bucket holding f_i . Thus, $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C$. If there is not, HeavyKeeper decays this bucket with a certain probability if and only if its counter is the first smallest among the counters of the d mapped buckets. Thus $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C$ or $A_j[h_j(f_i)](v+1).C = A_j[h_j(f_i)](v).C - 1$.

Therefore, $A_j[h_j(f_i)](v+1).C \leq n_i(v+1)$ holds in the minimum version. \square

C. Error Bound of the Basic Version of HeavyKeeper

Definition A.1. Given a small positive number ϵ , $Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\}$ ($n_i \geq \hat{n}_i$) represents the probability that the error of the estimated flow size $n_i - \hat{n}_i$ is larger than ϵN . If $Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} \leq \delta$, the algorithm is said to achieve (ϵ, δ) -counting.

(ϵ, δ) -counting is a metric to evaluate the error rate of the algorithm. Here HeavyKeeper is proved to achieve (ϵ, δ) -counting, showing that HeavyKeeper achieves a low error rate in estimating the sizes of top- k flows.

Theorem 5. Let's assume that there is no fingerprint collision and once the fingerprint of an elephant flow is inserted into its mapped bucket, it is held there all the time. Let's focus on one single array of HeavyKeeper. Given a small positive number ϵ , and an elephant flow f_i whose size is n_i is held at that bucket,

$$Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} \leq \frac{1}{\epsilon w n_i (b-1)} \quad (15)$$

where w is the width of each array, N the total number of packets, and b the exponential base, M the total number of different flows.

Proof. Let's focus on the j^{th} array. Flow f_i is correctly reported, so at the end, the fingerprint of flow f_i is held in the

$h_j(f_i)^{\text{th}}$ bucket of the j^{th} array. Let $I_{i,j,i'}$ be a binary random variable, defined as

$$I_{i,j,i'} = \begin{cases} 0 & (f_i = f_{i'}) \vee (h_j(f_i) \neq h_j(f_{i'})) \\ 1 & (f_i \neq f_{i'}) \wedge (h_j(f_i) = h_j(f_{i'})) \end{cases} \quad (16)$$

$I_{i,j,i'} = 1$ iff different flows f_i and $f_{i'}$ are held at the same bucket in the j^{th} array. We define random variable $X_{i,j}$ as:

$$X_{i,j} = \sum_{i'=1}^M I_{i,j,i'} n_{i'} \quad (17)$$

$X_{i,j}$ represents the sum of the sizes of the flows held at the same bucket as flow f_i , except for the size of flow f_i itself. Assume that for each incoming packet, if it belongs to flow f_i , the counter field is incremented by 1; if not, the counter field is decayed with a certain probability. We have

$$n_i - X_{i,j} \leq A_j[h_j(f_i)].C \leq n_i \quad (18)$$

Specifically, if all packets that do not belong to flow f_i decay the counter field, then $A_j[h_j(f_i)].C = n_i - X_{i,j}$. If those packets do not decay the counter field, then $A_j[h_j(f_i)].C = n_i$. Let's define another random variable $P_{i,j,l}$. Among the $X_{i,j}$ packets defined above, $P_{i,j,l}$ is defined as the probability that the l^{th} packet decays the counter field. Therefore,

$$A_j[h_j(f_i)].C = n_i - \sum_{l=1}^{X_{i,j}} P_{i,j,l} \quad (19)$$

Given a small positive number ϵ , the following formula based on the Markov inequality holds

$$\begin{aligned} & Pr\{A_j[h_j(f_i)].C \leq n_i - \epsilon N\} \\ &= Pr\{n_i - \sum_{l=1}^{X_{i,j}} P_{i,j,l} \leq n_i - \epsilon N\} \\ &= Pr\{\sum_{l=1}^{X_{i,j}} P_{i,j,l} \geq \epsilon N\} \leq \frac{E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})}{\epsilon N} \end{aligned} \quad (20)$$

Now let's focus on $E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$. Assume that all packets are uniformly distributed, for $\forall C$ satisfying $1 \leq C \leq n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$, we have the following formula:

$$Pr\{P_{i,j,l} = \frac{1}{b^C}\} = \frac{1}{A_j[h_j(f_i)].C} = \frac{1}{n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})} \quad (21)$$

Let β be $n_i - E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})$ for convenience. As a result,

$$\begin{aligned} E(\sum_{l=1}^{X_{i,j}} P_{i,j,l}) &= \sum_{l=1}^{E(X_{i,j})} E(P_{i,j,l}) \\ &= E(X_{i,j}) \sum_{c=1}^{\beta} \frac{1}{b^c} \frac{1}{\beta} = \frac{E(X_{i,j})}{\beta} \cdot \sum_{c=1}^{\beta} \frac{1}{b^c} \\ &= \frac{E(X_{i,j})}{\beta} \cdot \frac{\frac{1}{b}(1 - (\frac{1}{b})^{\beta})}{1 - \frac{1}{b}} \\ &\leq \frac{E(X_{i,j})}{n_i b} \cdot \frac{1 - (\frac{1}{b})^{n_i}}{1 - \frac{1}{b}} = \frac{E(X_{i,j})(1 - (\frac{1}{b})^{n_i})}{n_i(b-1)} \end{aligned} \quad (22)$$

Furthermore, for $E(X_{i,j})$, based on Equation 17,

$$E(X_{i,j}) = E\left(\sum_{i'=1}^M I_{i,j,i'} n_{i'}\right) \leq \sum_{i'=1}^M n_{i'} E(I_{i,j,i'}) = \frac{N}{w} \quad (23)$$

Therefore, based on Equation 22,

$$E\left(\sum_{l=1}^{X_{i,j}} P_{i,j,l}\right) \leq \frac{N(1 - (\frac{1}{b})^{n_i})}{wn_i(b-1)} \leq \frac{N}{wn_i(b-1)} \quad (24)$$

then

$$\begin{aligned} Pr\{A_j[h_j(f_i)].C \leq n_i - \epsilon N\} &\leq \frac{E(\sum_{l=1}^{X_{i,j}} P_{i,j,l})}{\epsilon N} \\ &\leq \frac{N}{\epsilon N wn_i(b-1)} = \frac{1}{\epsilon wn_i(b-1)} \end{aligned}$$

Note that for an elephant flow f_i , n_i is very large, and $(\frac{1}{b})^{n_i} \approx 0$. The estimated size of f_i is the maximum value of $A_j[h_j(f_i)].C$, so we have

$$Pr\{n_i - \hat{n}_i \geq \lceil \epsilon N \rceil\} \leq Pr\{\hat{n}_i \leq n_i - \epsilon N\} \leq \frac{1}{\epsilon wn_i(b-1)}$$

□

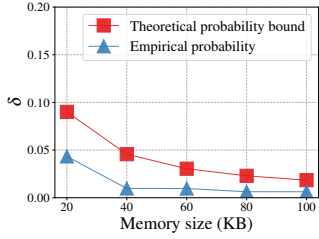


Fig. 35. Theoretical bound and empirical probability of HeavyKeeper ($\epsilon = 2^{-16}$).

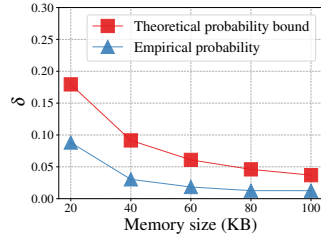


Fig. 36. Theoretical bound and empirical probability of HeavyKeeper ($\epsilon = 2^{-17}$).

To validate the correctness of this error bound, we conduct experiments on the dataset mentioned in Section VI-A. Here, we let $N = 10^7$, $\epsilon = 2^{-16}$ and 2^{-17} , and vary memory size from 20KB to 100KB. As shown in Figure 35 and Figure 36, the empirical probability of the basic version of HeavyKeeper is always lower than the theoretical probability bound, confirming the correctness of Theorem 5. Moreover, for the CSS algorithm, achieving such a (ϵ, δ) -counting requires at least $O(\epsilon^{-1})$ buckets (*i.e.*, $m = O(\epsilon^{-1})$), which requires a memory size much larger than 100KB. Therefore, HeavyKeeper is much more memory efficient than CSS.