

M4: A Framework for Per-Flow Quantile Estimation

Siyuan Dong*, Zhuochen Fan*, Tianyu Bai*, Tong Yang*, Hanyu Xue[†], Peiqing Chen*[‡], Yuhan Wu*

*School of Electronics Engineering and Computer Science, School of Computer Science, Peking University, China

[†]School of Mathematical Sciences, Peking University, China [‡]Department of Computer Science, University of Maryland, USA

Abstract—The field of quantile estimation has grown in importance due to its myriad practical applications. Recent research trends have evolved from estimating the quantile for a single data stream to developing data structures that can concurrently estimate quantiles for multiple sub-streams, also known as flows. This paper introduces a novel framework, M4, designed to estimate per-flow quantiles in data streams accurately. M4 is a versatile framework that can be integrated with a wide array of single-flow quantile estimation algorithms, thereby enabling them to perform per-flow estimation. The framework employs a sketch-based approach to provide a space-efficient method for recording and extracting distribution information. M4 incorporates two techniques: *MINIMUM* and *SUM*. The *MINIMUM* technique minimizes the noise on a flow from other flows caused by hash collisions, while the *SUM* technique efficiently categorizes flows based on their sizes and customizes treatment strategies accordingly. We demonstrate the application of M4 on three single-flow quantile estimation algorithms (DDSketch, *t*-digest, and ReqSketch), detailing the specific implementation of the *MINIMUM* and *SUM* techniques. We provide theoretical proof that M4 delivers high accuracy while utilizing limited memory. Additionally, we conduct extensive experiments to evaluate the performance of M4 regarding accuracy and speed. The experimental results indicate that across all three example algorithms, M4 significantly outperforms two comparison frameworks in terms of accuracy for per-flow quantile estimation while maintaining comparable speed.

I. INTRODUCTION

A. Background and Motivation

With the development of data stream processing, accurate, real-time extraction of required information from a large volume of high-speed data streams is attracting increasing attention [1]–[5]. Among the various types of information, quantile information, which requires distribution statistics for data streams, has become a focal point of numerous studies [6]–[12]. Recent research trends have evolved from estimating the quantile for a single data stream to developing data structures that can concurrently estimate quantiles for multiple sub-streams, also known as flows. In practical scenarios, many metrics necessitate per-flow granularity estimation of distribution, such as Latency [13]–[22], Inter-Arrival Time [23]–[25], Packet Size [26]–[29], and TTL (Time to Live) Value [30], [31]. Accurate estimation of per-flow distribution has wide-ranging practical applications and significant potential in distributed network scenarios, including improving the quality of service (QoS) for users [32]–[34], enhancing network

anomaly detection [35]–[37], and boosting the performance of Content Delivery Networks (CDNs) [38]. Consequently, the primary objective of this article is to perform quantile estimation for each individual flow in the data stream.

A data stream is a sequence of items, each represented as a *key-value* pair. Items sharing the same *key* compose a *flow*. The shared *key* serves as the *flow ID*¹. The *value* is the metric that needs processing. All items from different flows are intermixed in a data stream (e.g., $DS = \{\langle a, 3 \rangle, \langle a, 2 \rangle, \langle b, 5 \rangle, \langle d, 1 \rangle, \langle a, 4 \rangle, \dots\}$). This paper uses *quantile* to demonstrate per-flow *value* distribution. The items in a flow can be represented by a multiset $\mathcal{F} = \{\langle a, x_1 \rangle, \langle a, x_2 \rangle, \dots, \langle a, x_n \rangle\}$ of size n , where a is the *key*, x_i is the *value* and $x_1 \leq x_2 \leq \dots \leq x_n$. Given a percentage p ($0 \leq p \leq 1$), the p -quantile of *value* is x^* s.t. the percentage of $x_i \leq x^*$ in the multiset \mathcal{F} equals to p . With the above preliminaries, we give the problem definitions:

- **Per-Flow Quantile Estimation.**

```
SELECT key, p-quantile(value)
FROM DataStream
GROUP BY key
```

- **Single-Flow Quantile Estimation.**

```
SELECT p-quantile(value)
FROM DataStream
```

Accurately estimating the per-flow distribution is of wide practical usage and has many important potential applications in distributed scenarios. We provide three use cases as follows:

1) Improving of the quality of service (QoS) for online app users. In the digital era, online applications such as real-time video communications, online gaming, and streaming services have become integral to daily life [32]–[34]. These applications demand high-quality network services to ensure seamless user experience. Any degradation in network performance, particularly in terms of latency, can significantly affect the Quality of Service (QoS). This is especially true for applications requiring real-time interactions, such as remote control systems and remote sensing applications, where delays can compromise operational integrity and user experience. The challenge for network managers in maintaining optimal QoS lies in the ability to precisely identify and rectify latency

¹A flow ID is typically defined as a part of the five tuples: source IP address, destination IP address, source port, destination port, and protocol. This paper considers the number of items in a flow as the flow size, also referred to as the item frequency.

issues at a per-user granularity. By pinpointing the exact source and user with latency issues, network managers can implement targeted interventions to resolve these issues swiftly and efficiently.

2) More effective network anomaly detection. In this scenario, a flow refers to packets going through a network link with the same five-tuple ID. Network anomalies, such as congestion or the presence of malware, often manifest as abrupt increases in the latency of several flows [35]–[37] within a network link. These anomalies are critical to identify and mitigate, as they can severely impact network performance and security. Traditional single-flow quantile estimation approaches aggregate data across all flows, which can dilute the impact of anomalies present in a small subset of flows. This aggregation effect can lead to a failure in detecting subtle yet critical anomalies, allowing them to persist undetected and potentially cause extensive damage to the network infrastructure or compromise sensitive data. Our algorithm, by contrast, ensures that even minor deviations in latency are detected, thereby significantly enhancing the sensitivity of anomaly detection.

3) Cache performance optimization. A core challenge in this domain [39]–[41] is the diverse latency requirements and usage patterns across different flows, where a flow may represent a distinct user or a specific service. High latency in a flow often signals suboptimal data placement within the cache, necessitating strategic adjustments to either the data’s location or its storage modality to enhance access speed. Our approach effectively evaluates cache strategy impacts on different flows, identifying when optimizations reduces average latency but inadvertently disadvantage high-priority users. It also identifies strategies that benefit latency-insensitive services at the cost of sensitive ones, allowing the development of cache strategies that enhance system efficiency without compromising critical service performance.

Numerous studies have significantly advanced the field of single-flow quantile estimation [15]–[18], [42]–[45]. Approximate algorithms—often referred to as sketches—have predominantly excelled in scenarios requiring low memory overhead and rapid processing, with only minimal accuracy trade-offs. However, a critical limitation inherent in these approaches is their inability to discern among multiple flow IDs. They are designed either to estimate quantiles for individual keys in isolation or to aggregate across all data stream values without key differentiation. This results in a binary choice: a focused but isolated single-flow estimation or a comprehensive yet undifferentiated analysis across all flows. It fails to address the nuanced requirements of modern networked systems, where identifying and analyzing per-flow metrics is crucial for optimizing performance and detecting anomalies.

The evolution of applications and the diversity of their requirements necessitate a more granular approach to quantile estimation—one that can accurately measure and adapt to the unique characteristics of each flow. To bridge this gap, there are two strategic paths: designing a brand new algorithm tailored for per-flow estimation [46], [47] or a

versatile framework capable of enabling existing single-flow algorithms to handle per-flow queries. We advocate for the latter, recognizing its potential for broad applicability and design simplicity. This approach leverages the strengths of existing algorithms to cater to scenarios with varying priorities like throughput, accuracy, and memory efficiency.

By opting for a framework that transforms single-flow algorithms into their per-flow counterparts, we present a solution that is both adaptable and scalable. This framework not only retains the inherent advantages of the original algorithms but also expands their utility to support per-flow estimation, thereby addressing a critical need in network management and analysis. In doing so, our approach provides a comprehensive toolset for network administrators and researchers, enabling them to tackle a wide array of challenges with unprecedented precision and flexibility.

There are several challenges when designing such a framework. **(1) Algorithmic Compatibility.** Different single-flow algorithms may have unique characteristics, optimizations, and assumptions that may not directly translate to a per-flow context. The framework must provide a flexible architecture that can adapt various single-flow algorithms to per-flow requirements without compromising their inherent advantages. **(2) Scalability.** One of the foremost challenges is ensuring that the framework scales efficiently with the number of flows. Single-flow algorithms are typically optimized for performance with a single data stream. Extending these to accommodate multiple, potentially thousands or millions of flows, can introduce significant computational and memory overhead. The framework must efficiently manage resources to maintain high performance and accuracy across all flows. **(3) Data Skew and Flow Variability.** In real-world networks, some flows may be more active or larger than others, leading to data skew. The framework needs to handle such variability, ensuring that large or high-volume flows do not overshadow smaller ones.

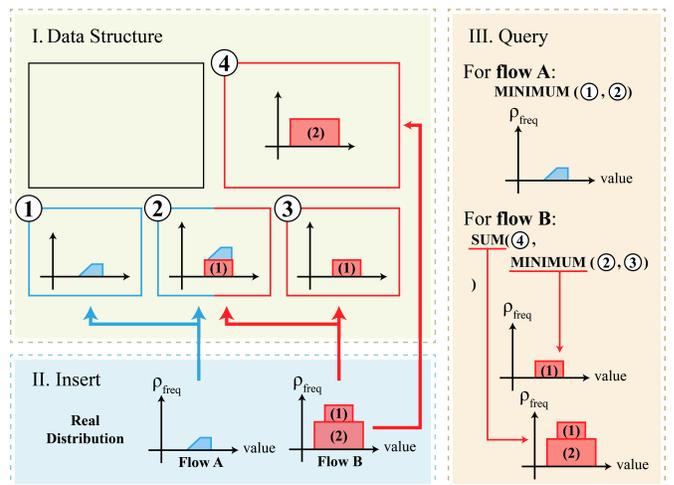


Fig. 1: **Illustration of M4.** Arrows pointing from II to I represent the hashing operation. Flow A is represented by blue. Flow B is represented by red.

B. Our Solution and Contributions

To achieve our design goal, we propose a novel framework named MINIMUM-SUM (M4). *M4 is a framework that can be applied to an extensive range of single-flow quantile estimation algorithms, enabling them to perform per-flow quantile estimation.* For simplicity, we refer to the single-flow algorithm on which we employ M4 as *META*. As depicted in Figure 1, M4 uses limited memory to construct several layers of buckets. Each bucket contains a *META* to record distribution. Every *META* treats all incoming items identically. We use hash functions to map flows to buckets for recording. A single flow can be mapped to multiple buckets. If a hash collision occurs, the distribution of the collided flows will sum up in the bucket. Each bucket has a load capacity. The insertion of a flow starts at a lower layer. When the buckets in the lower layer overflow, we insert the subsequent items into the upper layers. Buckets in higher layers have larger capacity and finer granularity. M4 comprises two techniques: *MINIMUM* and *SUM*. As long as the *META* can generate the *value* distribution for a single flow, we can use *MINIMUM* and *SUM* to transform it into an efficient per-flow quantile estimation algorithm.

MINIMUM resolves hash collisions by randomly selecting several buckets for each flow in each layer using multi-hashes and extracting the real distribution from these selected buckets. To address hash collisions, a straightforward solution is to use a hash table of buckets, each containing a *META* to record the *value* distribution of a flow. However, this approach requires recording IDs and executing complex operations (typically dependent on the number of flows) to locate another available bucket during a hash collision, which is both memory and time consuming. A more efficient solution is *Memory Sharing*. We use w hash functions to map a flow to w buckets of *META* for recording, providing us with w records for every flow. Due to hash collisions, more than one flow may be mapped into one bucket. Thus, every flow record may contain some noise from other flows, and we need to compare and analyze the distributions given by all the records to restore the real distribution. For instance, as shown in Figure 1, a small flow A and a large flow B are inserted into M4². Flow A can be contained in the first layer, while flow B is segmented into parts I and II³ due to its large size. Flow A and B encounter a hash collision at bucket ② in the first layer, so their distribution information is mixed in the bucket. We need to use bucket ① and ② to restore the real distribution of flow A and use bucket ② and ③ to restore the real distribution of part I of flow B. Because noise at one *value* point only increases the density at that point, by selecting the lowest density at each *value* point, the estimated *value* distribution bears the slightest noise possible. As we take the intersection of distributions, we call it the *MINIMUM* technique. We can see in Figure 1 that the

²There is a predefined threshold categorizing flows according to their sizes. Flows with sizes below that threshold are called small flows, otherwise large flows.

³We cut flow B into top and bottom halves just for simplicity.

MINIMUM technique allows us to restore the real distribution of flow A and part I of flow B. It should be noted that if *META* can give an error-free result when estimating single-flow distribution, our *MINIMUM* is optimal and can guarantee a unilateral error of over-estimation when estimating per-flow distribution.

SUM categorizes flows based on their sizes through segmentation. It uses multiple layers to segment every flow into multiple parts and tailor the treatment in each layer. Then it aggregates these parts to generate the full distribution. There are two reasons why we need flow categorization. First, the flow size distribution in a data stream is usually highly skewed. For example, in CAIDA [48] dataset, about 40% flows only contain ≤ 3 items. It would be a colossal waste to allocate equal resources to a small flow and a large one. Second, large flows get more attention in practical applications⁴. Hence, the larger a flow is, the more resources should be allocated to achieve a fair or better accuracy than small flows. Thus, we want to categorize flows according to their sizes efficiently. However, we do not know if a flow is large or small in advance. To solve this, M4 has a layered structure, where lower layers are for coarse-grained recording of small flows and higher layers are for fine-grained recording of large flows. Each flow is considered a small flow in the beginning. When the recorded flow size exceeds the capacity in lower layers, we know it is a large flow, so we insert the subsequent items into higher layers. In this way, the distribution information of a large flow is scattered at multiple layers. We need to aggregate all parts together to get the entire distribution. As shown in Figure 1, we need to use bucket ④ and (②*MINIMUM*③) to restore the real distribution of flow B. Because the distribution in two layers corresponds to two disjoint parts of flow B, we should sum up the density at each *value* point to construct the overall distribution information. Therefore, we call it the *SUM* technique. We can see in Figure 1 that the *SUM* technique allows us to restore the real distribution of flow B.

We apply our M4 to three *METAs* (DDSketch [15], t -digest [16], [18], and modified ReqSketch [17], [42]). The three *METAs* each have their emphasis. DDSketch allows us to focus on the tail *value* distribution and bounds the relative error of quantile estimation to a constant at different percentages. t -digest allows us to tailor the relative accuracy of quantile estimation at different percentages. ReqSketch provides a relative guarantee on the error of rank estimation. We design the *MINIMUM* and *SUM* techniques for them according to their features.

We devise two comparison frameworks (see Section V-A2 and Section V-A3) for better comparison. We perform extensive experiments to evaluate our performance regarding accuracy and speed. Experiment results indicate that *M4 is per-flow friendly and accurate*. For tiny flows, maximum *value* estimation is on average 90.6% error-free, while comparison frameworks offer almost no error-free estimates. For larger

⁴Large flows tend to represent critical or high-priority data. By paying more attention to large flows, administrators can optimize the delivery of important information and enhance overall system performance.

flows, the Average Logarithm Error (ALE) of M4 reach $2.26 \times$ lower than comparison frameworks. *M4 is memory-efficient*. It only needs 6MB to handle 27M items. We also provide theoretical proof that M4 delivers high accuracy while utilizing limited memory.

Key contributions:

- We introduce M4, the first general framework that can be applied to a wide range of single-flow quantile estimation algorithms to accomplish per-flow quantile estimation, filling a gap in the research field.
- We propose the *MINIMUM* and *SUM* techniques. Together, they reduce the error from hash collisions and allow us to tailor treatment strategies for flows of different sizes.
- We apply M4 to DDSketch, *t*-digest, and modified ReqSketch and implement them on a CPU platform. Compared to two comparison frameworks, M4 achieves significantly better accuracy with a comparable speed across all three algorithms. All codes are available on GitHub [49].

II. RELATED WORK

A. Sketch

A sketch is a type of probabilistic data structure designed to process data with small and controllable errors. One of the most classic sketches is CM Sketch [50], designed for estimating item frequency. CM Sketch consists of d arrays, each array $A_i (1 \leq i \leq d)$ has w counters and is associated with a hash function $h_i(\cdot)$. When an incoming item e is inserted, we increase the counter $A_i[h_i(e) \% w]$ by 1 for all $i \in \{1, 2, \dots, d\}$. To query the item e , CM Sketch reports the minimum counter among all the d mapped counters determined by hash functions. Other classic sketches include Flajolet-Martin (FM) Sketch [51], CU Sketch [52], Count Sketch [53], CSM Sketch [54] and CMM Sketch [55].

B. Single-Flow Quantile Estimation

Quantile estimation approximates the distribution of metrics, essential for analyzing large or streaming datasets. While previous studies have advanced single-flow quantile estimation [15]–[18], [42]–[45], they fall short in multi-flow scenarios typical in data streams from diverse sources. Our framework, M4, extends these methods to estimate distributions per flow. We illustrate its application using three algorithms: DDSketch, *t*-digest, and ReqSketch, each with unique advantages. DDSketch allows us to focus on the tail *value* distribution and bounds the relative error of quantile estimation to a constant at different percentages. *t*-digest allows us to tailor the relative accuracy of quantile estimation at different percentages. ReqSketch provides a relative guarantee on the error of rank estimation.

1) DDSketch:

DDSketch [15] is designed for estimating value distributions by partitioning the value range into segments, each monitored by a counter for values within the segment. The segment boundaries are determined by $\gamma := (1 + \alpha)/(1 - \alpha)$, with each segment's counter, C_i , tallying values x in the range $\gamma^{i-1} < x \leq \gamma^i$. The insertion of a value x is indexed by

$\lceil \log_\gamma(x) \rceil$. DDSketch approximates values in a segment by $\hat{x} = \frac{2\gamma^i}{\gamma+1}$, maintaining a relative error within α . The trade-off between the range coverage and accuracy is governed by α : a higher α extends the range but reduces accuracy. To estimate the value at a certain percentile $p \in [0, 1]$, we sum counters up to the relevant segment i and use $\hat{x} = \frac{2\gamma^i}{\gamma+1}$ as the quantile estimate.

2) *t*-digest:

The *t*-digest [16], [18] algorithm clusters real-valued samples to approximate value distributions, grouping items by similarity. Each cluster records the mean *value* and total count of its items. Items are added to the closest cluster, updating its statistics. *t*-digest controls cluster sizes to balance precision and memory use, adjusting cluster counts through a scale function k and a compression parameter δ . The function k ensures uniform cluster weight growth, allowing more clusters where data is denser. For quantile queries, weights are summed until the target cluster is identified, assuming uniform distribution within clusters to estimate the queried value.

t-digest [16], [18] is designed to estimate *value* distributions by clustering real-valued samples. *t*-digest uses clusters to group items with near *values*. Each cluster contains an *average* cell recording the mean *value* of absorbed items, and a *weight* cell recording the total number of absorbed items. Each incoming item is assigned to the cluster with the nearest average *value*, after which the average and weight cells of that cluster are updated. The key idea of *t*-digest is to confine the weight of each cluster to an appropriate level, being small enough to record the distribution accurately, while large enough to avoid unacceptable memory costs. Accurate confinement is achieved by constantly monitoring all clusters' weights and keeping them at the same level. There are a non-decreasing *scale function* $k : [0, 1] \rightarrow \mathbb{R}$ describing the weight restriction and a *compression parameter* δ bounding the number of clusters used. We define w_i as cluster C_i 's *weight* value, and N as $\sum_i w_i$. Each w_i must satisfy: $k(\frac{w_{<i}+w_i}{N}) - k(\frac{w_{<i}}{N}) \leq \frac{1}{\delta}$. As a result, *t*-digest allocates more clusters to the segment of *value* with more items. Besides, we can tailor the relative accuracy of quantile estimation at different percentages by changing the *scale function* k . To query for the *value* at percentage $p \in [0, 1]$, we accumulate cluster weights until finding the cluster that p falls in. Deeming that items are uniformly distributed in each cluster, we can get the estimated *value* according to the position of p in that cluster.

3) ReqSketch:

ReqSketch [17], [42] employs multiple levels of *compactors* to store item *values*, utilizing $O(\log N)$ *compactors* for a flow size N , each with a buffer of similar size. Items enter at level 0 and, upon a compactor's capacity being reached, a sorted even-sized subset is compacted and half its elements are elevated to the next level, preserving total item weight due to the weighting scheme where items at level h are assigned a weight of 2^h . This mechanism, known as compaction, ensures the integrity of distribution estimates. To estimate a value at a

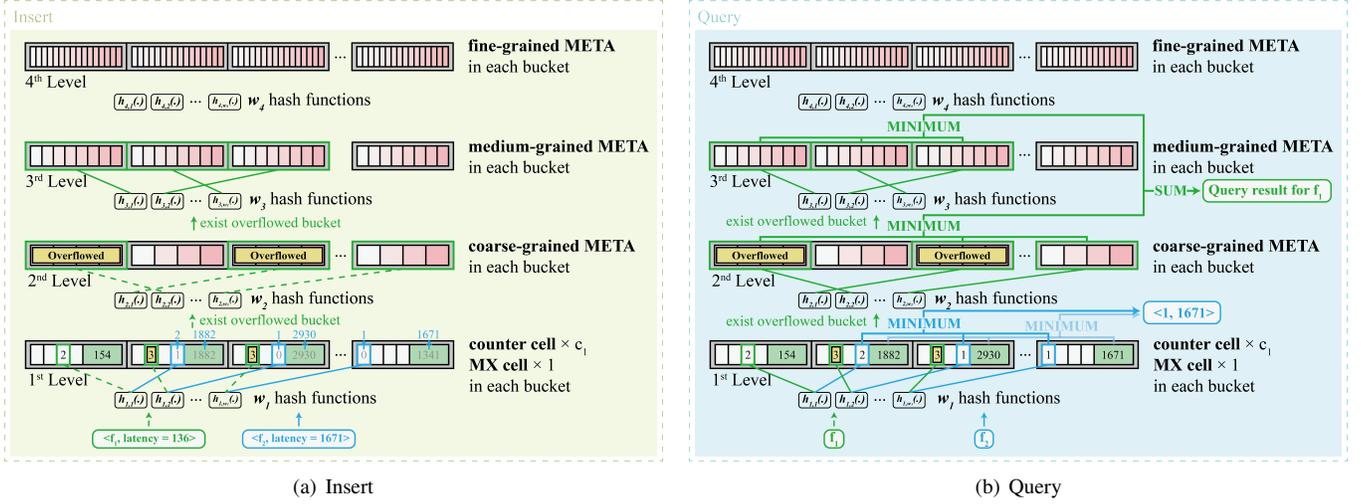


Fig. 2: Framework structure of M4.

certain percentile p , item weights are cumulatively tallied from the smallest value until reaching p , with the corresponding item's value serving as the estimate.

III. M4 DESIGN

A. Problem Statement

DEFINITION 1. Data Stream. A data stream is a series of items appearing in sequence. Each item e_i is a key-value pair. The key serves as an ID, while the value represents the metric we aim to process. An example of a data stream is $DS = \{\langle a, 3 \rangle, \langle a, 2 \rangle, \langle b, 5 \rangle, \langle d, 1 \rangle, \langle a, 4 \rangle, \dots\}$.

DEFINITION 2. Flow. Items sharing the same key compose a flow, and the shared key is their flow ID. The number of items in a flow is the flow size, also called the item frequency. An example of a flow is $\mathcal{F} = \{\langle a, 3 \rangle, \langle a, 2 \rangle, \langle a, 4 \rangle, \dots\}$.

DEFINITION 3. Quantile. Given a numerical multiset $S = \{x_1, x_2, \dots, x_n\}$ of size n , where $x_1 \leq x_2 \leq \dots \leq x_n$, and a percentage p ($0 \leq p \leq 1$), the p -quantile of multiset S is defined as $x_{\lfloor p(n-1) \rfloor + 1}$.

Per-Flow Quantile Estimation. Given an arbitrary flow f in a data stream of key-value pairs and a percentage p , we need to estimate the p -quantile of value in f . To express it in SQL:

```
CREATE TABLE DataStream (
  key int,
  value int
)
/* Insert items into DataStream. */
SELECT key, p-quantile(value)
FROM DataStream
GROUP BY key
```

B. Framework Description

This section describes the structure and operations of M4. Frequently used notations are outlined in Table I.

TABLE I: Notations.

Symbol	Meaning
e	an item in the data stream
f	the key (flow ID) of a certain item
v	the value of a certain item
L_i	the i^{th} level of bucket array
b_i	the number of buckets in L_i
w_i	the number of hash functions associated with L_i
$h_{i,j}(\cdot)$	the j^{th} hash function in L_i
l_i	the capacity parameter for buckets in L_i
c_i	the granularity parameter for buckets in L_i

1) Framework Structure:

In Figure 2, M4 is presented as a four-tiered bucket array, chosen for its balance between precision and efficiency, as further discussed in Section V-B. Each level L_i comprises b_i buckets and employs w_i hash functions ($h_{i,1}(\cdot)$ to $h_{i,w_i}(\cdot)$). L_1 captures the size and maximum value of small flows, L_2 the value distribution of moderate flows, and L_3 and L_4 the distribution of large flows. We classify tiny, medium, and huge flows as having sizes in $[1, 3)$, $[3, 255)$, and $[255, +\infty)$, respectively. The algorithm M4, when integrated with META (DDSketch, t -digest, or modified ReqSketch), is referred to as M4-META.

Each L_1 bucket in any M4-META contains c_1 counters and an MX cell for flow size and maximum value, respectively, with a counter cell of l_1 bits. Overflows occur when a cell's count maxes out. L_1 is designed for scenarios where small flows are less critical.

For levels L_i where $i \geq 2$, each bucket holds a META, with uniform capacity and granularity within the level but varying across levels to accommodate the significance and size of the flows. A bucket overflows once its META reaches capacity. The details for different METAs on L_i ($i \geq 2$) are as follows: **DDSketch**. Referencing Section II-B1, DDSketches in levels L_i ($i \geq 2$) consist of c_i counter cells ($c_2 \leq c_3 \leq c_4$), each tracking the frequency of values within distinct segments. The bit length of counters in L_i is l_i ($l_2 \leq l_3 \leq l_4$). A bucket in L_i overflows when any counter cell's frequency hits its maximum l_i -bit value.

***t*-digest.** As per Section II-B2, a *t*-digest in levels $L_i (i \geq 2)$ includes c_i clusters ($c_2 \leq c_3 \leq c_4$), with each cluster holding an *average* and a *weight* cell for the *mean* value and item count, respectively. The weight cell length in L_i is l_i ($l_2 \leq l_3 \leq l_4$). Overflow occurs when a bucket’s weight cell frequency in L_i maxes out its l_i -bit capacity.

mReqSketch. Mentioned in Section II-B3, mReqSketches in $L_i (i \geq 2)$ comprise l_i compactors ($l_2 \leq l_3 \leq l_4$) with c_i cells each ($c_2 \leq c_3 \leq c_4$) for value storage. The maximum weight for an mReqSketch in L_i is $(2^{l_i} - 1) \times c_i$. Overflow is determined when a bucket’s recorded frequency in L_i reaches this limit.

2) *Framework Insertion Operation:*

To insert an item $e = \langle f, v \rangle$ into M4-META, we target the lowest non-overflowed level, denoted as L_{top} . First, e is mapped to w_1 buckets in L_1 using the index $\lfloor \frac{h_{1,j}(f)\%c_1}{c_1} \rfloor$ for $j \in 1, 2, \dots, w_1$. The $(h_{1,j}(f)\%c_1)^{th}$ counter cell in each mapped bucket is incremented, and the *MX* cell is updated to $\max MX, v$, unless overflow occurs, prompting an attempt to insert e into L_2 .

For levels L_i where $i \geq 2$, e is mapped to w_i buckets using $h_{i,j}(f)\%b_i$. If no overflow occurs in the mapped buckets, e is inserted into the *META* of each, concluding the insertion.

3) *Framework Query Operation:*

Querying a flow with ID f involves mapping it to corresponding buckets across levels, starting from L_1 up to L_{top} , the highest non-overflowed level for f . Unlike insertion, querying aggregates results from all relevant levels up to L_{top} .

At each level L_i , we obtain w_i records for f , which may be affected by hash collisions. To mitigate this, we employ the *MINIMUM* technique to derive the least polluted distribution by selecting the minimum values from counter and *MX* cells for L_1 , representing the size and maximum value of f . For levels L_i where $i \geq 2$, the approach adjusts based on *META*’s structure, detailed in Section III-C.

The output for f depends on *top*. For $top = 1$, the result is based solely on L_1 data using the *MINIMUM* method. For $top \geq 2$, it is a *SUM*-merged aggregation from L_2 to L_{top} , with specifics on the *SUM* technique in Section III-C.

4) *Example:*

Our example uses parameters $\langle c_1 = 4, l_1 = 2, w_1 = 3 \rangle$, $\langle w_2 = 3 \rangle$, $\langle w_3 = 3 \rangle$, and $\langle w_4 = 3 \rangle$.

In the insertion operation illustrated in Figure 2(a), item $e_1 = \langle f_1, v = 136 \rangle$ is first mapped to L_1 ’s three buckets, but due to counter overflows, it’s redirected and successfully inserted into L_3 . For item $e_2 = \langle f_2, v = 1671 \rangle$, mapping to L_1 reveals no overflow, so counters are incremented, and 1341 is updated to 1671.

Figure 2(b) shows the query operation. For f_1 , overflow at L_1 and L_2 leads to $top = 3$, with results merged using the *SUM* and *MINIMUM* techniques from L_2 and L_3 . For f_2 , L_1 provides the flow size as the minimum counter value (1) and the maximum value as the smallest *MX* cell (1671).

C. *MINIMUM & SUM*

In this section, we expound on applying *MINIMUM* and *SUM* to an arbitrary *META* and illustrate the workflow on three

example *META*s, DDSketch, *t*-digest, and mReqSketch. First, we present the distribution stored in the *META* as histograms (Section III-C2). Subsequently, we perform *MINIMUM* and *SUM* operations on the histograms (Section III-C3).

1) *Rationale:*

MINIMUM. The *MINIMUM* technique mitigates hash collision effects without tracking IDs, facilitating $O(1)$ time complexity for both insertion and query operations. Recognizing that accurate distribution estimation equates to frequency estimation at each value point, this method leverages the fact that lower densities in flow buckets, resulting from hash collisions, provide a more accurate density at any given value point. Thus, selecting the minimum density from all mapped buckets yields the most reliable value distribution estimate.

SUM. The *SUM* technique efficiently categorizes flows based on their sizes and tailors treatment strategies accordingly, maximizing the overall accuracy. To achieve so, we use multiple layers to divide a flow’s distribution information into various fractions. Since the information in these fractions is disjointed, we need to sum up the density at each *value* point to construct the overall distribution, akin to piecing together a jigsaw puzzle.

2) *From META to Histogram:*

Histograms are a widely used method for representing distributions. Consequently, every *META* can transform the distribution information stored with its data structure into a histogram. In this subsection, we illustrate how DDSketch, *t*-digest, and mReqSketch are transformed into histograms.

DDSketch. As discussed in Section II-B1, DDSketch divides the entire range of *value* into fixed segments, each tracked by a counter cell that records the number of *values* that fall into that segment. If we index each segment by $i \in \mathbb{Z}$, then the counter C_i records the number of *value* x that falls between $V_{i-1} = \gamma^{i-1} < x \leq \gamma^i = V_i$.

The process of transforming a DDSketch into a histogram is illustrated on the left side of Figure 3. The range of each segment on the horizontal axis is determined by V_i ($i \in \{0, 1, 2, 3\}$), and the frequency of each segment is the corresponding C_i .

***t*-digest.** As discussed in Section II-B2, *t*-digest uses clusters to group items with near *values*. Each cluster contains an *average* cell V_i recording the mean *value* of absorbed items, and a *weight* cell C_i recording the total number of absorbed items. Each incoming item is assigned to the cluster with the nearest average *value*, after which the V_i and C_i of that cluster are updated. Besides, *t*-digest records the minimum *value* V_{min} and the maximum *value* V_{max} .

The process of transforming a *t*-digest into a histogram is illustrated in the middle of Figure 3. The range of each segment on the horizontal axis is determined by V_{min} , V_{max} , and V_i ($i \in \{1, 2, 3\}$). Since V_i is an average *value*, we divide C_i into two halves and distribute them to adjacent segments.

mReqSketch. As discussed in Section II-B3, ReqSketch consists of several levels of *compactor* serving as buffers. Each level comprises multiple cells storing the *value* V_i of items. Each cell in level i carries a weight $C_i = 2^i$ ($i \in \{0, 1, 2, \dots\}$).

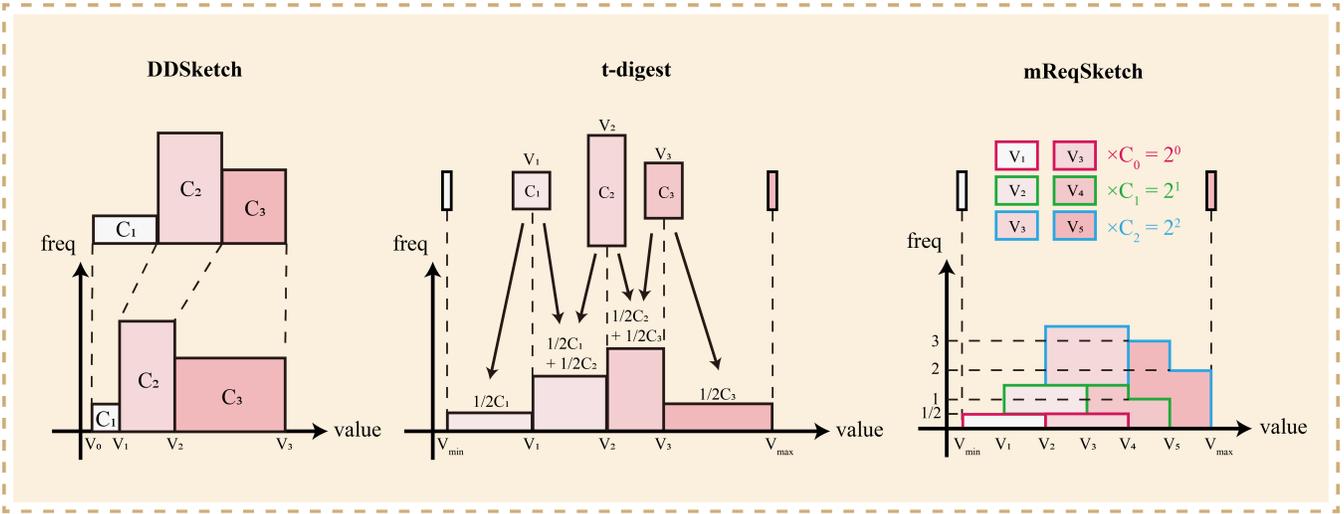


Fig. 3: From *META* to Histogram.

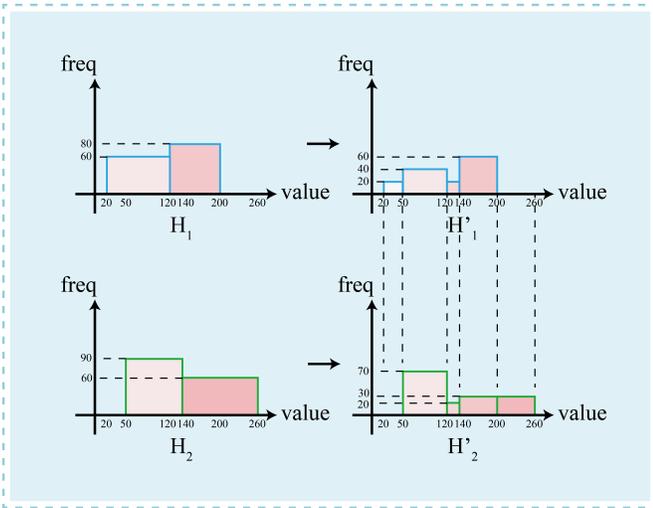


Fig. 4: Segment Alignment.

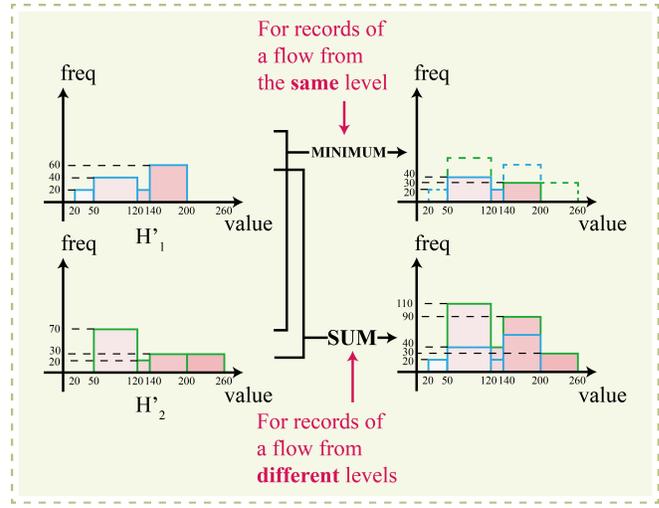


Fig. 5: *MINIMUM* & *SUM* on aligned Histograms.

Besides, ReqSketch also records the minimum V_{min} and the maximum V_{max} .

The original design of *compaction operation* in ReqSketch is memory-intensive and slow, making it unsuitable for estimating per-flow *value* distribution in data streams. We introduce minor modifications to the original design while maintaining its quintessence and rename it mReqSketch. According to the new design, the incoming items are always inserted into level 0. Whenever a level h becomes full, we sort the items in level h , remove them from this level, and randomly select half of the items (either odd or even indexed) to be inserted into level $h + 1$.

The process of transforming a mReqSketch into a histogram is illustrated on the right side of Figure 3. The range of each segment on the horizontal axis is determined by V_{min} , V_{max} , and V_i ($i \in \{1, 2, 3, 4, 5\}$). Since V_i is a randomly selected *value*, we divide C_i into two halves and distribute them to adjacent segments.

3) *MINIMUM* & *SUM* on Histogram:

After transforming *META* into histograms, we can perform the *MINIMUM* and *SUM* operations. These operations necessitate a prerequisite known as *Segment Alignment*. As illustrated in Figure 4, suppose we have two histograms H_1 and H_2 , which require alignment. We first need to obtain the union boundary set on the horizontal axis. The boundary set in H_1 is $S_1 = \{20, 120, 200\}$, and that of H_2 is $S_2 = \{50, 140, 260\}$. Thus, the union boundary set is $U = \{20, 50, 120, 140, 200, 260\}$. Next, for both histograms, we scatter the frequency recorded in each segment determined by S_i into the new segments determined by U , following a uniform distribution. The rationale here is that focusing on a minimal interval allows us to approximate any distribution with a uniform one, mirroring the central concept in calculus.

After *Segment Alignment*, we can operate on the frequency in each segment of H'_1 and H'_2 . As shown in Figure 5, if H'_1 and H'_2 originate from the same level of the same flow, we

will need to *MINIMUM*-merge them. We select the minimum frequency in each segment to construct the *MINIMUM*-merged distribution. If H'_1 and H'_2 originate from different levels of the same flow, we will need to *SUM*-merge them. We add the frequencies in each segment to construct the *SUM*-merged distribution.

To query for the *value* at percentage $p \in [0, 1]$, we accumulate segment frequencies until finding the segment that p falls in. Then we calculate a quantile according to the uniform distribution and report it as the estimation result.

D. Discussion

1) Using histogram as an intermediate step:

A critical aspect of developing a universally applicable framework is ensuring the adaptability of core operations, such as *MINIMUM* and *SUM*, to accommodate a wide range of single-flow algorithms. The inherent challenge lies in the diversity of distribution outputs that quantile estimation algorithms can produce. Our solution is normalizing these distributions into histograms, which serve as a standardized intermediary representation. It simplifies the implementation of the *MINIMUM* and *SUM* operations by reducing them to operations on histogram bins, thereby enhancing the framework's efficiency and scalability. This design choice underscores the flexibility of the M4 framework, allowing it to seamlessly incorporate a vast array of quantile estimation algorithms without necessitating extensive customization or reconfiguration. This adaptability is crucial for a general framework aimed at broad applicability across diverse network environments and applications.

2) Data Aging:

One potential area for improvement in our framework is enhancing the data aging process to better align with real-time data distribution changes. Currently, the mechanism for phasing out outdated data is clearing the data structure for every time window. It may not be sufficiently prompt, creating a gap between the stored data distribution and the actual current distribution. This discrepancy becomes particularly problematic when data distribution shifts rapidly, potentially compromising the method's effectiveness. Improving the data aging process to more accurately reflect immediate distribution changes is essential for maintaining M4's accuracy in dynamic data environments.

3) Flow size distribution:

Our methodology, including the underlying data structures and operations, is specifically optimized for data streams characterized by a long-tailed distribution of flow sizes, a common phenomenon in real-world applications where the principle of *the few governing the many* often applies. This long-tailed distribution pattern ensures that a small fraction of flows carries a significant portion of data.

It is important to note, however, that in environments where the flow size distribution is uniform, our approach may not be the most efficient. We acknowledge that this specificity may limit the universal applicability of our approach but we

believe that its optimized performance in its intended context represents a significant contribution to the field.

IV. MATHEMATICAL ANALYSIS

There are two sources of error for M4: (1) The accuracy of the recording algorithm in each bucket and (2) hash collision. The error introduced by (1) is due to inaccuracies in *META*. We will present an error bound for mReqSketch in this section as we modified the original design of ReqSketch. As for (2), which is hash collision related, we will present an error bound for M4-DDSketch. The analysis is conducted at one specific level at once, assuming that there are n buckets and m flows arriving at the level we are examining, with each flow mapped to w buckets.

A. Analysis of M4-DDSketch

We begin by defining some frequently used notations. $freq$ represents the total number of items at a level. $freq_i$ represents the number of items in flow i at this level. $freq_i^T$ represents the number of items in flow i within a segment T at this level. $freq^T$ represents the total number of items within the segment T at this level.

1) Huge and Medium Flows:

Theorem 1. Let \hat{freq}_i denote the estimation of $freq_i$. Then

$$P(\hat{freq}_i > freq_i + \epsilon) < \left(\frac{freq}{n\epsilon}\right)^w \quad (1)$$

Proof. Let $\hat{freq}_{i,k}$ ($k \in \{1, 2, \dots, w\}$) denote the k^{th} record of $freq_i$, independent from each other. Since $\hat{freq}_i = \min\{\hat{freq}_{i,k}\}$, we obtain that

$$P(\hat{freq}_i > freq_i + \epsilon) = P^w(\hat{freq}_{i,k} > freq_i + \epsilon) \quad (2)$$

Now we analyze the situation where $w = 1$. Let A_j denote the event that flow j ($j \neq i$) is mapped to the same bucket as flow i . Then we have $\hat{freq}_{i,1} = freq_i + \sum_{j \neq i} freq_j 1_{A_j}$ and $P(A_j) = \frac{1}{n}$. Hence, $E(\hat{freq}_{i,1} - freq_i) = \sum_{j \neq i} freq_j \frac{1}{n} < \frac{freq}{n}$. Because $\hat{freq}_{i,1} - freq_i \geq 0$, we obtain

$$P(\hat{freq}_{i,1} - freq_i > \epsilon) \leq \frac{E(\hat{freq}_{i,1} - freq_i)}{\epsilon} < \frac{freq}{n\epsilon} \quad (3)$$

Therefore,

$$P(\hat{freq}_i > freq_i + \epsilon) = P^w(\hat{freq}_{i,1} - freq_i > \epsilon) < \left(\frac{freq}{n\epsilon}\right)^w \quad (4)$$

□

Theorem 2. Let \hat{freq}_i^T denote the estimation of $freq_i^T$. Then

$$P(\hat{freq}_i^T > freq_i^T + \epsilon) < \left(\frac{freq^T}{n\epsilon}\right)^w \quad (5)$$

Proof. Let $\hat{freq}_{i,k}^T$ ($k \in \{1, 2, \dots, w\}$) denote the k^{th} record of $freq_i^T$, independent from each other. Since $\hat{freq}_i^T = \min\{\hat{freq}_{i,k}^T\}$, we obtain

$$P(\hat{freq}_i^T > freq_i^T + \epsilon) = P^w(\hat{freq}_{i,k}^T > freq_i^T + \epsilon) \quad (6)$$

Now we analyze the situation where $w = 1$. Let A_j denote the event that flow j ($j \neq i$) is mapped to the same bucket as flow i . Then we have $\hat{freq}_{i,1}^T = freq_i^T + \sum_{j \neq i} freq_j^T 1_{A_j}$ and $P(A_j) = \frac{1}{n}$. Hence, $E(\hat{freq}_{i,1}^T - freq_i^T) = \sum_{j \neq i} freq_j^T \frac{1}{n} < \frac{freq^T}{n}$. Because $\hat{freq}_{i,1}^T - freq_i^T \geq 0$, we obtain

$$P(\hat{freq}_{i,1}^T - freq_i^T > \epsilon) \leq \frac{E(\hat{freq}_{i,1}^T - freq_i^T)}{\epsilon} < \frac{freq^T}{n\epsilon} \quad (7)$$

Therefore,

$$\begin{aligned} P(\hat{freq}_i^T > freq_i^T + \epsilon) &= P^w(\hat{freq}_{i,1}^T - freq_i^T > \epsilon) \\ &< \left(\frac{freq^T}{n\epsilon}\right)^w \end{aligned} \quad (8)$$

□

Theorem 3. Let \hat{t}_p denote the estimated quantile of percentage p . Then

$$P(|\hat{t}_p - t_p| < \alpha t_p) \geq 1 - (1 - e^{-\frac{m}{n}})^w \quad (9)$$

Proof. The probability of hash collision happening in all mapped buckets of a flow is $P_C = [1 - (\frac{n-1}{n})^m]^w \approx (1 - e^{-\frac{m}{n}})^w$. So the probability that there is at least one bucket where no hash collision occurs is $1 - P_C = 1 - (1 - e^{-\frac{m}{n}})^w$. In this case, the error is bounded by $|\hat{t}_p - t_p| < \alpha t_p$, as proved in DDSketch. Therefore, $P(|\hat{t}_p - t_p| < \alpha t_p) \geq 1 - (1 - e^{-\frac{m}{n}})^w$ □

Hash collisions may have a significant impact on quantile t_p . Unless strong assumptions are made on the *value* distributions of flows, giving an error bound of \hat{t}_p is impossible when hash collisions happen in all mapped buckets.

Comparison with prior work. We choose to compare M4-DDSketch with SketchPolymer [47], a per-flow quantile estimation algorithm which has a similar algorithm principle with M4-DDSketch. The main difference between M4-DDSketch and SketchPolymer is that M4-DDSketch uses multiple layers for huge and medium flows, while SketchPolymer only uses one layer after the filtration of tiny flows. When the memory is limited, our algorithm is usually more accurate for huge flows. The reason is that by using less bits storing medium flows, we can allocate more memory for huge flows to avoid hash collisions.

First, let us consider M4-DDSketch. We divide the flows to medium and huge flows. The DDSketch for medium and huge flows has the same α^5 , but we use fewer bits for the counters of medium flows. We denote the memory cost of one bucket for medium and huge flows by b_1 and b_2 , respectively. We set $w = 1$ and the amount of buckets for medium and huge flows to be n_1 and n_2 . We denote the total memory by C . Then we have

$$n_1 b_1 + n_2 b_2 = C \quad (10)$$

We can choose the proportion of n_1 and n_2 to make the expected value of hash collision frequency (denoted by $E(HC)$) the same for medium and huge flows. For a huge flow i , the error of its frequency (sum of medium and huge parts) $freq_i$ has the expectation value:

$$E(\hat{freq}_i - freq_i) = E(HC)F, \quad (11)$$

where F is the total number of items. We denote the number of medium flows by X_1 and the number of huge flows by X_2 , then we have

$$\begin{aligned} X_1 + X_2 &= E(HC)n_1 \\ X_2 &= E(HC)n_2 \end{aligned} \quad (12)$$

Combining equation (10), (11) and (12), we get

$$E(\hat{freq}_i - freq_i) = \frac{F}{C}(X_1 b_1 + X_2 b_1 + X_2 b_2) \quad (13)$$

Therefore, we have the error bound

$$\begin{aligned} P(\hat{freq}_i - freq_i > \epsilon) &\leq \frac{E(\hat{freq}_i - freq_i)}{\epsilon} \\ &= \frac{F}{C\epsilon}(X_1 b_1 + X_2 b_1 + X_2 b_2) \end{aligned} \quad (14)$$

Next, let us consider SketchPolymer. We have total memory = C , bucket size = b_2 (all the buckets should be of the size for huge flows), X_1 medium flows and X_2 huge flows. The corresponding error bound is

$$P(\hat{freq}_i' - freq_i > \epsilon) \leq \frac{F}{C\epsilon}(X_1 b_2 + X_2 b_2) \quad (15)$$

The difference of the error bound is

$$\frac{F}{C\epsilon}(X_2 b_1 - X_1(b_2 - b_1)) \quad (16)$$

In real situations, X_1 is much larger than X_2 . So the error bound of M4-DDSketch is smaller than that of SketchPolymer.

2) *Tiny Flows:*

Theorem 4. Let $x = \frac{wm}{n}$, then the probability of getting a wrong maximum value is $(\frac{e^{-x} + x - 1}{x})^w$.

Proof. Suppose that a flow f is mapped to w buckets a_1, a_2, \dots, a_w . The probability that there are k_i other flows in bucket a_i ($i \in \{1, 2, \dots, w\}$) is

$$\begin{aligned} P_{\{k_i\}} &= \frac{\binom{wm-w}{k_1} \binom{wm-w-k_1}{k_2} \dots \binom{wm-w-k_1-k_2-\dots-k_{w-1}}{k_w}}{n^{wm-w}} \\ &\cdot (n-w)^{wm-k_1-k_2-\dots-k_w-w}. \end{aligned} \quad (17)$$

The equation above is based on the assumption that the w hash values of a flow are independent. Due to the factor $(n-w)^{-k_1-k_2-\dots-k_w}$, the probability of k_i being large is low. So we can assume that $k_i \ll m$ and arrive at the approximation $P_{\{k_i\}} \approx \prod_i \binom{wm-w}{k_i} (n-w)^{-k_i} (1 - \frac{w}{n})^{wm-w} \approx \prod_i \binom{wm}{k_i} n^{-k_i} e^{-\frac{w^2 m}{n}}$. In this situation, the probability of

⁵For the definition of α , please refer to Section II-B1

obtaining an incorrect maximum *value* is $\prod_i \frac{k_i}{k_i+1}$. Hence, the probability of obtaining an incorrect maximum *value* is

$$\begin{aligned}
& \left[\prod_i \sum_{k_i=1}^{wm} \binom{wm}{k_i} n^{-k_i} \frac{k_i}{k_i+1} \right] e^{-\frac{w^2 m}{n}} \\
&= \left[\sum_{k=1}^{wm} \binom{wm}{k} n^{-k} \frac{k}{k+1} \right]^w e^{-\frac{w^2 m}{n}} \\
&= \left[\frac{n - (1 + \frac{1}{n})^{wm} (-wm + n)}{wm + 1} e^{-\frac{wm}{n}} \right]^w \\
&\stackrel{x=\frac{wm}{n}}{\approx} \left(\frac{e^{-x} + x - 1}{x} \right)^w
\end{aligned} \tag{18}$$

□

B. Analysis of *mReqSketch*

In this section, we conduct the error analysis for *mReqSketch*. $R(y)$ represents the count of *values* that is less than or equal to *value* y across all items recorded at a level. $R_i(y)$ represents the count of *value* that is less than or equal to *value* y in flow i at this level.

Error bound of *mReqSketch*. Consider the following setting. There are M compactors C_0, C_2, \dots, C_{M-1} in the *mReqSketch*. The buffer size in each compactor is $b = 2^a$. The number of items in this *mReqSketch* is at maximum capacity $N = 2^a(2^M - 1) \approx 2^{a+M}$. Let p denote the real value of the fraction of *values* less than y . Consider the estimation for $R(y)$.

Theorem 5. Let $\hat{R}(y)$ denote the estimation of $R(y)$. Then

$$P(|\hat{R}(y) - R(y)| > \epsilon N) < 2e^{-\frac{4b^2}{1-(1-2p)^b} \epsilon^2} \tag{19}$$

Each time we conduct the *compaction operation*, we operate on b *values*. Among these, the probability that the number of *value* less than y is odd is given by

$$\begin{aligned}
P &= \frac{\sum_{i=1}^b \left[\binom{b}{i} p^i (1-p)^{b-i} \right] - \sum_{i=1}^b \left[\binom{b}{i} (-p)^i (1-p)^{b-i} \right]}{2} \\
&= \frac{1 - (1-2p)^b}{2}
\end{aligned} \tag{20}$$

We must perform the *compaction operation* on C_0 for 2^M times. If the selected items have odd indices, the probability of having an error of $+1$ on $R(y)$ is P , and the probability of error-free is $(1-P)$. If items with even indices are selected, the probability of having an error of -1 on $R(y)$ is P , and the probability of error-free is $(1-P)$. Therefore, the expected error is 0 each time, and the error variance is $P(1-P)$. The overall expected error is 0, and the overall error variance is $2^M P(1-P)$.

We need to perform the *compaction operation* on C_1 for 2^{M-1} times. If items with odd indices are selected, the probability of having an error of $+2$ on $R(y)$ is P , and the probability of error-free is $(1-P)$. If items with even indices are selected, the probability of having an error of -2 on $R(y)$ is P , and the probability of error-free is $(1-P)$. Therefore,

TABLE II: Default parameter settings.

Algorithm	c			l		
	c_2	c_3	c_4	l_2	l_3	l_4
M4-DDSketch	20	20	35	8	16	32
M4- <i>t</i> -digest	4	8	16	8	16	32
M4-mReqSketch	2	2	4	8	16	32
Strawman-DDSketch	35			32		
Strawman- <i>t</i> -digest	16			32		
Strawman-mReqSketch	4			32		
CuckooFilter-DDSketch	68			32		
CuckooFilter- <i>t</i> -digest	32			32		
CuckooFilter-mReqSketch	8			32		

the expected error is 0 each time, and the error variance is $4P(1-P)$. The overall expected error is 0, and the overall error variance is $2^{M+1}P(1-P)$.

Performing the above analysis for all the compactors, we find that the overall expected error in $C_i (i \in \{0, 1, 2, \dots, M-1\})$ is 0. The overall variance of the error in C_i is $2^{M+i}P(1-P)$. Therefore, the variance of the error across the whole *mReqSketch* is

$$\begin{aligned}
\sigma^2 &= P(1-P)(2^M + 2^{M+1} + \dots + 2^{2M-1}) \\
&= P(1-P)2^M(2^M - 1) \\
&\approx \frac{1 - (1-2p)^{2b}}{4} 2^{2M} \\
&\approx \frac{1 - (1-2p)^b}{4} \frac{N^2}{b^2}
\end{aligned} \tag{21}$$

Note that the variance from C_i increases with i , so the Lindeberg-Feller condition is not satisfied. We cannot apply the central limit theorem. However, we can employ the sub-Gaussian distribution estimation and find

$$P(|\hat{R}(y) - R(y)| > \epsilon N) < 2e^{-\left(\frac{\epsilon N}{\sigma}\right)^2} = 2e^{-\frac{4b^2}{1-(1-2p)^b} \epsilon^2} \tag{22}$$

V. EXPERIMENTAL RESULTS

A. Experiment Setup

1) Implementation:

We implemented M4 and all related *META* data structures (DDSketch, *t*-digest, and *mReqSketch*) in C++. The hash functions were devised using the 32-bit Bob Hash (sourced from an open-source website [56]), each with different initial seeds. All the experiments were executed on an 18-core CPU server (Intel i9-10980XE) with 128GB memory and 24.75MB L3 cache. Each experiment was repeated ten times to compute an average result.

2) Straw-man Solution:

To compare effectively, we developed a straw-man solution using a Dleft-like approach for each *META* (DDSketch, *t*-digest, and *mReqSketch*), featuring three bucket arrays for an optimal balance between accuracy and speed (details in Section V-B). Each bucket records a flow ID *key* and the distribution of *values* using a *META*. We use three distinct hash functions, $h'_1(\cdot)$, $h'_2(\cdot)$, and $h'_3(\cdot)$, for the arrays. Additionally, a global *META* aggregates the value distribution for all flows, serving as a fallback for queries on unrecorded flows.

Insertion. When a new item $e = \langle f, v \rangle$ arrives, it's first added to the global *META*. We then try to insert e to the

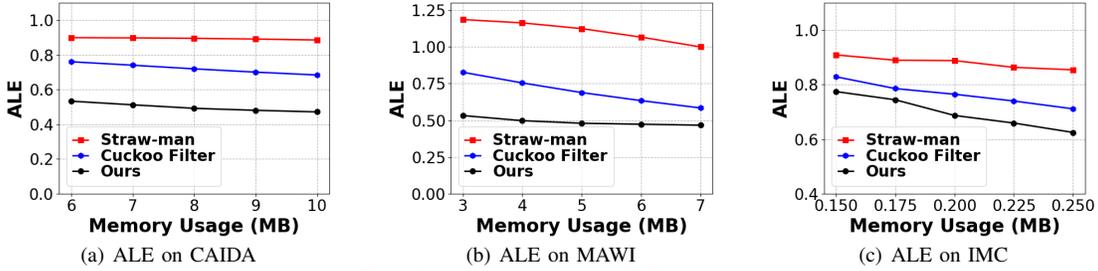


Fig. 6: Accuracy of M4-DDSketch.

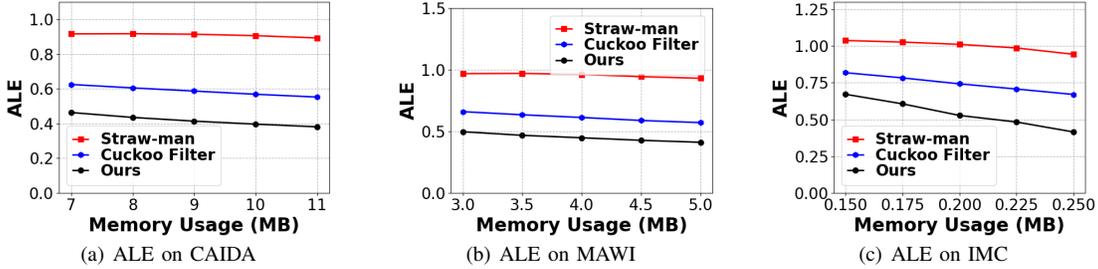


Fig. 7: Accuracy of M4-t-digest.

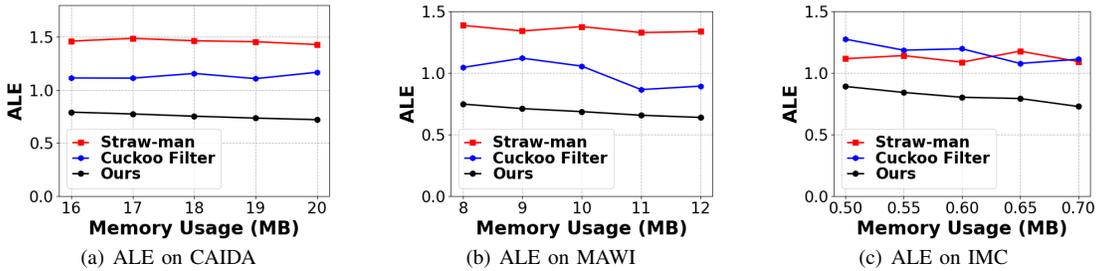


Fig. 8: Accuracy of M4-mReqSketch.

$h'_1(f)^{th}$ bucket of the first array. If this bucket is free or already contains a flow with $key = f$, e is added to this bucket's META, setting key to f . Failing that, we move to the $h'_2(f)^{th}$ bucket of the second array, and if necessary, to the $h'_3(f)^{th}$ bucket of the third array. Should all arrays reject e , it is discarded.

Query. To find a flow f , we check the $h'_1(f)^{th}$, $h'_2(f)^{th}$, and $h'_3(f)^{th}$ bucket of the respective arrays for a $key = f$. A match returns the query from that bucket's META; otherwise, the global META provides the result.

3) Cuckoo Filter:

We devised another comparison framework based on Cuckoo Filter [57]. Cuckoo Filter is an efficient hash table implementation based on cuckoo hashing [58], which can achieve both high utilization and compactness. It records the fingerprint instead of the flow ID to improve space efficiency.

The structure employs a table with buckets and three hash functions, $h_1(\cdot)$, $h_2(\cdot)$, and $h_f(\cdot)$, where each bucket records a fingerprint fp^c and the distribution of values using a META. For an item $e = \langle f, v \rangle$, we determine $fp = h_f(f)$ and map e to the $[h_1(fp)]^{th}$ and $[h_2(fp)]^{th}$ buckets. We then try to locate one of these two buckets where $fp^c = fp$ and insert v to the META in that bucket. If no matching bucket is found but an empty bucket exists, we insert e to it by setting its $fp^c = fp$ and inserting v to the META. If both buckets

are full, then one of the two flows in them will be evicted to its alternate bucket (because each flow has two mapped buckets). This random eviction continues until an empty bucket is found or a preset maximum number of attempts, $MAX_NUMBER_OF_TURNS$, is reached. If that happens, we discard this item. Like what we do in the straw-man solution, a global META is maintained as a fallback for queries. We chose $MAX_NUMBER_OF_TURNS = 8$ and 32-bit fp to minimize collisions, as increasing $MAX_NUMBER_OF_TURNS$ beyond this point does not significantly enhance accuracy but does reduce throughput.

4) Datasets:

- 1) *CAIDA Dataset.* This dataset comprises streams of anonymized IP items collected from high-speed monitors by CAIDA in 2018 [48]. We use the trace with a monitoring interval of 60s. Each item consists of a 5-tuple (13 bytes). There are around 27M items and 1.3M flows in this dataset.
- 2) *MAWI Dataset.* This dataset contains real traffic trace data maintained by the MAWI Working Group [59]. Similar to CAIDA, each item in the dataset is a 5-tuple. There are around 9M items and 13K flows in the MAWI dataset.
- 3) *IMC Dataset.* This dataset comes from one of the data centers studied in [60]. Each item also consists of a 5-tuple. There are around 14M items and 5K flows in this

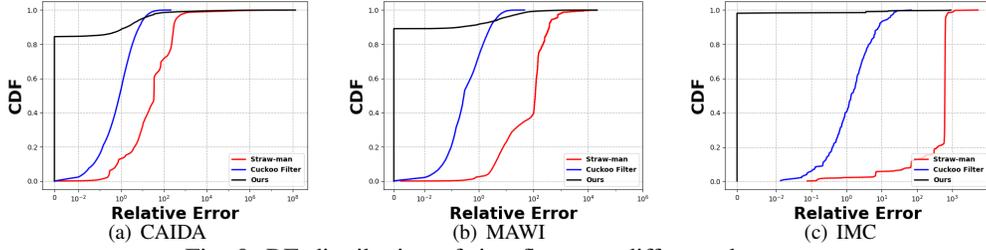


Fig. 9: RE distribution of tiny flows on different datasets.

dataset.

5) Metrics:

- 1) *ALE (Average Logarithm Error)*. We employ *ALE* to evaluate the accuracy of quantile estimation for huge and medium flows. Since the order magnitude of latency may vary significantly, it is unreasonable to measure the error by absolute value alone. We define *ALE* as $\frac{1}{|\Psi|} \sum_{f_i \in \Psi} |\log_2 t_i - \log_2 \hat{t}_i| = \frac{1}{|\Psi|} \sum_{f_i \in \Psi} |\log_2 \frac{t_i}{\hat{t}_i}|$, where Ψ represents the set of all huge and medium flows in the data stream. t_i and \hat{t}_i denote the real and estimated quantile at a given percentage p .
- 2) *RE (Relative Error)*. We employ *RE* to evaluate the accuracy of quantile estimation for tiny flows. We define *RE* as $\frac{|\hat{x}_{max} - x_{max}|}{x_{max}}$, where x_{max} and \hat{x}_{max} are the real and estimated maximum value in a flow.
- 3) *Throughput (Insertion and Query)*. We use million operations (insert an item or query a flow) per second (Mops) to measure the throughput.

6) Default settings:

In our experiment, we set $p = 0.5$ as the default setting⁶. For M4, the memory ratio of L_1, L_2, L_3, L_4 is 3%, 60%, 35%, 2%, respectively. Each arriving item at every level is mapped to 3 buckets ($w_i = 3, i \in \{1, 2, 3, 4\}$). Each bucket in L_1 has 4 counter cells ($c_1 = 4$) and 1 MX cell, and each counter cell consists of 2 bits ($l_1 = 2$). The parameter settings for DDSketch, *t*-digest, and mReqSketch on different frameworks are shown in Table II.

B. Experiments on Parameter Setting

We conduct experiments on (1) the choice of the number of levels for M4 and the straw-man solution, (2) the choice of the number of hash functions (w) in each layer, and (3) the choice of p on the performance. Due to limited space, we move this part of experiments to our appendix [49].

C. Experiments of Huge and Medium Flows on Accuracy

Experiments on M4-DDSketch. As shown in Figure 6, the experimental results show that the ALE of M4-DDSketch is significantly lower than that of the comparison frameworks. Specifically on the three real-world datasets, the ALEs of M4-DDSketch are on average $1.80\times$, $2.26\times$, and $1.27\times$ lower than those of the straw-man solution, and $1.45\times$, $1.42\times$ and $1.10\times$ lower than Cuckoo Filter.

⁶Experiment results under different settings of p are similar. For more details, please refer to Section V-B

Experiments on M4-*t*-digest. As shown in Figure 7, the experimental results demonstrate that the ALE of M4-*t*-digest is significantly lower than that of the comparison frameworks. Specifically on the three real-world datasets, the ALEs of M4-*t*-digest are on average $2.19\times$, $2.13\times$ and $1.90\times$ lower than those of the straw-man solution, and $1.41\times$, $1.36\times$ and $1.40\times$ lower than Cuckoo Filter.

Experiments on M4-mReqSketch. As shown in Figure 8, the experimental results demonstrate that the ALE of M4-mReqSketch is significantly lower than that of the comparison frameworks. Specifically on the three real-world datasets, the ALEs of M4-mReqSketch are on average $1.94\times$, $1.97\times$, and $1.39\times$ lower than those of the straw-man solution, and $1.50\times$, $1.45\times$ and $1.44\times$ lower than Cuckoo Filter.

Analysis. The accuracy advantage in Figure 6-8 is established by making better use of resources. First, the separation of medium and huge flows brought by the *SUM* technique allows for allocating fewer bits for medium flows. Furthermore, the multi-layer structure prevents huge and medium flows from contaminating each other. Otherwise, once a medium flow collides with a huge one, its distribution will be utterly covered up by the huge flow because of its size. Second, the *MINIMUM* technique improves the algorithm's robustness against hash collisions.

D. Experiments of Tiny Flows on Accuracy

As shown in Figure 9, the experimental results demonstrate that for most tiny flows, the maximum value can be well estimated by M4, which significantly outperforms the two comparison frameworks. Specifically on the three real-world datasets, M4 attains an error-free (i.e., $RE = 0$) rate of 84.5%, 89.1%, and 98.2%, while the comparison frameworks offer almost no error-free estimates.

Tiny flows are too small to well-define a distribution. Using *METAs* to record tiny flows would be inaccurate and memory-consuming. Hence, only recording the maximum value gives us significantly better results than two comparison frameworks.

E. Experiments on Speed

We conduct experiments on the speed of M4 and two comparison frameworks on different datasets. Due to limited space, we move this part of experiments to the appendix [49]. The experiments show that the throughput of M4 in insertion and query is slightly lower than that of the comparison frameworks. This is because our solution contains more layers. Besides, *MINIMUM* and *SUM* operations take extra time.

VI. CONCLUSION

This paper introduces the M4 framework designed to enable per-flow quantile estimation using single-flow estimation algorithms. The key techniques of M4 are *MINIMUM*, employed for minimization of the noise caused by hash collisions, and *SUM*, employed for efficient flow categorization based on their sizes and customized treatment strategies. The experimental results indicate that M4 outperforms two comparison frameworks in estimating the *value* distribution of huge, medium, and tiny flows. For huge and medium flows, the ALE is $1.89\times$ lower than the straw-man solution, and $1.49\times$ lower than Cuckoo Filter. For tiny flows, the maximum *value* estimation attains an error-free rate of 98.0%, which is a stark improvement over the virtually nonexistent error-free estimates offered by the comparison frameworks. We have made our code publicly available on GitHub [49] to facilitate further research and application in this field.

ACKNOWLEDGMENT

We would like to thank the anonymous reviewers for their thoughtful feedback. This work was supported in part by the National Key R&D Program of China (No. 2022YFB2901504), in part by the National Natural Science Foundation of China (NSFC) (No. U20A20179, 62372009), and in part by the China Postdoctoral Science Foundation (No. 2023TQ0010, GZC20230055).

REFERENCES

- [1] D. Nguyen, G. Memik, S.O. Memik, and A. Choudhary. Real-time feature extraction for high speed networks. In *Proc. FPL*, pages 438–443, 2005.
- [2] Albert Bifet. Mining big data in real time. *informatica*, 37(1), 2013.
- [3] Eduardo Viegas, Altair Santin, Alysson Bessani, and Nuno Neves. Bigflow: Real-time and reliable anomaly-based intrusion detection for high-speed networks. *Future Generation Computer Systems*, 93:473–485, 2019.
- [4] M Mazhar Rathore, Hojae Son, Awais Ahmad, Anand Paul, and Gwanggil Jeon. Real-time big data stream processing using gpu with spark over hadoop ecosystem. *International Journal of Parallel Programming*, 46:630–646, 2018.
- [5] Zhuochen Fan, Ruixin Wang, Yalun Cai, Ruwen Zhang, Tong Yang, Yuhan Wu, Bin Cui, and Steve Uhlig. Onesketch: A generic and accurate sketch for data streams. *IEEE Transactions on Knowledge and Data Engineering*, 35(12):12887–12901, 2023.
- [6] Tyler Akidau, Alex Balikov, Kaya Bekiroğlu, Slava Chernyak, Josh Haberman, Reuven Lax, Sam McVeety, Daniel Mills, Paul Nordstrom, and Sam Whittle. Millwheel: Fault-tolerant stream processing at internet scale. *Proc. VLDB Endow.*, 6(11):1033–1044, 2013.
- [7] Jun Gao, Chang Zhou, Jiashuai Zhou, and Jeffrey Xu Yu. Continuous pattern detection over billion-edge graph using distributed framework. In *Proc. ICDE*, pages 556–567, 2014.
- [8] Jörn Kuhlenskamp, Markus Klems, and Oliver Röss. Benchmarking scalability and elasticity of distributed database systems. *Proc. VLDB Endow.*, 7(12):1219–1230, 2014.
- [9] Yuxiang Zeng, Yongxin Tong, and Lei Chen. Last-mile delivery made practical: An efficient route planning framework with theoretical guarantees. *Proc. VLDB Endow.*, 13(3):320–333, 2019.
- [10] Cong Yue, Zhongle Xie, Meihui Zhang, Gang Chen, Beng Chin Ooi, Sheng Wang, and Xiaokui Xiao. Analysis of indexing structures for immutable data. In *Proc. SIGMOD*, pages 925–935, 2020.
- [11] Jiayi Wang, Chengliang Chai, Jiabin Liu, and Guoliang Li. Face: A normalizing flow based cardinality estimator. *Proc. VLDB Endow.*, 15(1):72–84, 2021.
- [12] Bo Wang, Rongqiang Chen, and Lu Tang. Easyquantile: Efficient quantile tracking in the data plane. In *Proc. APNet*, pages 123–129, 2023.
- [13] Ying Zhang, Wenjie Zhang, Jian Pei, Xuemin Lin, Qianlu Lin, and Aiping Li. Consensus-based ranking of multivalued objects: A generalized borda count approach. *IEEE Transactions on Knowledge and Data Engineering*, 26(1):83–96, 2012.
- [14] Zubair Shah, Abdun Naser Mahmood, Zahir Tari, and Albert Y Zomaya. A technique for efficient query estimation over distributed data streams. *IEEE Transactions on Parallel and Distributed Systems*, 28(10):2770–2783, 2017.
- [15] Charles Masson, Jee E Rim, and Homin K Lee. Ddsksetch: A fast and fully-mergeable quantile sketch with relative-error guarantees. *Proc. VLDB Endow.*, 12(12):2195–2205, 2019.
- [16] Ted Dunning. The size of a t-digest. *arXiv preprint arXiv:1903.09921*, 2019.
- [17] Graham Cormode, Abhinav Mishra, Joseph Ross, and Pavel Vesely. Theory meets practice at the median: a worst case comparison of relative error quantile algorithms. In *Proc. KDD*, pages 2722–2731, 2021.
- [18] Ted Dunning. The t-digest: Efficient estimates of distributions. *Software Impacts*, 7:100049, 2021.
- [19] Michal Szymaniak, David Presotto, Guillaume Pierre, and Maarten van Steen. Practical large-scale latency estimation. *Computer Networks*, 52(7):1343–1364, 2008.
- [20] TS Eugene Ng and Hui Zhang. Predicting internet network distance with coordinates-based approaches. In *Proc. INFOCOM*, volume 1, pages 170–179, 2002.
- [21] Rui Zhu, Bang Liu, Di Niu, Zongpeng Li, and Hong Vicky Zhao. Network latency estimation for personal devices: A matrix completion approach. *IEEE/ACM Transactions on Networking*, 25(2):724–737, 2017.
- [22] Yongjun Liao, Wei Du, Pierre Geurts, and Guy Leduc. Dmfsgd: A decentralized matrix factorization algorithm for network distance prediction. *IEEE/ACM Transactions on Networking*, 21(5):1511–1524, 2013.
- [23] Gábor Horváth, Peter Buchholz, and Miklós Telek. A map fitting approach with independent approximation of the inter-arrival time distribution and the lag correlation. In *Proc. QEST*, pages 124–133, 2005.
- [24] Dimitris J Bertsimas and Garrett Van Ryzin. Stochastic and dynamic vehicle routing with general demand and interarrival time distributions. *Advances in Applied Probability*, 25(4):947–978, 1993.
- [25] Blake McShane, Moshe Adrian, Eric T Bradlow, and Peter S Fader. Count models based on weibull interarrival times. *Journal of Business & Economic Statistics*, 26(3):369–378, 2008.
- [26] Rishi Sinha, Christos Papadopoulos, and John Heidemann. Internet packet size distributions: Some observations. *USC/Information Sciences Institute, Tech. Rep. ISI-TR-2007-643*, pages 1536–1276, 2007.
- [27] Ping Du and Shunji Abe. Detecting dos attacks using packet size distribution. In *Proc. BIMNICS*, pages 93–96, 2007.
- [28] Chun-Nan Lu, Chun-Ying Huang, Ying-Dar Lin, and Yuan-Cheng Lai. Session level flow classification by packet size distribution and session grouping. *Computer Networks*, 56(1):260–272, 2012.
- [29] DJ Parish, K Bharadia, A Larkum, IW Phillips, and MA Oliver. Using packet size distributions to identify real-time networked applications. *IEE Proceedings-Communications*, 150(4):221–227, 2003.
- [30] Nicaise Choungmo Fofack, Philippe Nain, Giovanni Neglia, and Don Towsley. Performance evaluation of hierarchical ttl-based cache networks. *Computer Networks*, 65:212–231, 2014.
- [31] Edith Cohen, Eran Halperin, and Haim Kaplan. Performance aspects of distributed caches using ttl-based consistency. *Theoretical Computer Science*, 331(1):73–96, 2005.
- [32] Biao Wang, Ge Chen, Luoyi Fu, Li Song, and Xinning Wang. Drimux: Dynamic rumor influence minimization with user experience in social networks. *IEEE Transactions on Knowledge and Data Engineering*, 29(10):2168–2181, 2017.
- [33] Devinder Kaur, Gagangeet Singh Aujla, Neeraj Kumar, Albert Y Zomaya, Charith Perera, and Rajiv Ranjan. Tensor-based big data management scheme for dimensionality reduction problem in smart grid systems: Sdn perspective. *IEEE Transactions on Knowledge and Data Engineering*, 30(10):1985–1998, 2018.
- [34] Jing Li, Weifa Liang, Wenzheng Xu, Zichuan Xu, and Jin Zhao. Maximizing the quality of user experience of using services in edge

- computing for delay-sensitive iot applications. In *Proc. MSWiM*, pages 113–121, 2020.
- [35] Izzat Alsmadi and Dianxiang Xu. Security of software defined networks: A survey. *Computers Security*, 53:79–108, 2015.
- [36] Yixin Chen, Jianing Pei, and Defang Li. Detpro: a high-efficiency and low-latency system against ddos attacks in sdn based on decision tree. In *Proc. ICC*, pages 1–6, 2019.
- [37] Michael Kurth, Ben Gras, Dennis Andriess, Cristiano Giuffrida, Herbert Bos, and Kaveh Razavi. Netcat: Practical cache attacks from the network. In *Proc. SP*, pages 20–38, 2020.
- [38] Muhammad Shahzad and Alex X Liu. Noise can help: Accurate and efficient per-flow latency measurement without packet probing and time stamping. *ACM SIGMETRICS Performance Evaluation Review*, 42(1):207–219, 2014.
- [39] Avriella Floratou, Nimrod Megiddo, Navneet Potti, Fatma Özcan, Uday Kale, and Jan Schmitz-Hermes. Adaptive caching in big sql using the hdfs cache. In *Proc. SoCC*, pages 321–333, 2016.
- [40] Youtao Zhang and Rajiv Gupta. Enabling partial-cache line prefetching through data compression. *High-Performance Computing: Paradigm and Infrastructure*, pages 183–201, 2005.
- [41] Peiqing Chen, Dong Chen, Lingxiao Zheng, Jizhou Li, and Tong Yang. Out of many we are one: Measuring item batch with clock-sketch. In *Proc. SIGMOD*, pages 261–273, 2021.
- [42] Graham Cormode, Zohar Karnin, Edo Liberty, Justin Thaler, and Pavel Vesely. Relative error streaming quantiles. In *Proc. PODS*, pages 96–108, 2021.
- [43] Michael Greenwald and Sanjeev Khanna. Space-efficient online computation of quantile summaries. *ACM SIGMOD Record*, 30(2):58–66, 2001.
- [44] Zohar Karnin, Kevin Lang, and Edo Liberty. Optimal quantile approximation in streams. In *Proc. FOCS*, pages 71–78, 2016.
- [45] Rana Shahout, Roy Friedman, and Ran Ben Basat. Squad: Combining sketching and sampling is better than either for per-item quantile estimation. *arXiv preprint arXiv:2201.01958*, 2022.
- [46] Jintao He, Jiaqi Zhu, and Qun Huang. Histsketch: A compact data structure for accurate per-key distribution monitoring. In *Proc. ICDE*, pages 2008–2021, 2023.
- [47] Jiarui Guo, Yisen Hong, Yuhan Wu, Yunfei Liu, Tong Yang, and Bin Cui. Sketchpolymer: Estimate per-item tail quantile using one sketch. In *Proc. KDD*, pages 590–601, 2023.
- [48] The CAIDA Anonymized Internet Traces. <https://www.caida.org/catalog/datasets/overview/>.
- [49] The appendix and the source codes of ours along with other algorithms. <https://github.com/M4Framework/M4>.
- [50] Graham Cormode and Shan Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- [51] Philippe Flajolet and G Nigel Martin. Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences*, 31(2):182–209, 1985.
- [52] Cristian Estan and George Varghese. New directions in traffic measurement and accounting: Focusing on the elephants, ignoring the mice. *ACM Transactions on Computer Systems (TOCS)*, 21(3):270–313, 2003.
- [53] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Proc. ICALP*, pages 693–703, 2002.
- [54] Tao Li, Shigang Chen, and Yibei Ling. Per-flow traffic measurement through randomized counter sharing. *IEEE/ACM Transactions on Networking*, 20(5):1622–1634, 2012.
- [55] Fan Deng and Davood Rafiei. New estimation algorithms for streaming data: Count-min can do more. *Webdocs. Cs. Ualberta. Ca*, 2007.
- [56] Bob jenkins’ hash function web page, paper published in dr dobb’s journal. <http://burtleburtle.net/bob/hash/evahash.html>.
- [57] Bin Fan, Dave G Andersen, Michael Kaminsky, and Michael D Mitzenmacher. Cuckoo filter: Practically better than bloom. In *Proc. CoNEXT*, pages 75–88, 2014.
- [58] Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. *Journal of Algorithms*, 51(2):122–144, 2004.
- [59] MAWI Working Group Traffic Archive. <http://mawi.wide.ad.jp/mawi/>.
- [60] Theophilus Benson, Aditya Akella, and David A Maltz. Network traffic characteristics of data centers in the wild. In *Proc. IMC*, pages 267–280, 2010.